

GENERATIVE TRAFFIC SIMULATIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Controlling urban traffic is an emerging challenge for modern cities. In this work, we present a cascaded AI system that integrates a deep neural network with a large language model to generate and simulate both routine and edge-case traffic scenarios. The proposed system analyzes real-time congestion patterns and generates optimized traffic signal plans, adjusting signal timing based on vehicle flow direction and phase-level congestion to minimize delays and enhance throughput. By combining real-world traffic data with synthetically generated scenes, the approach enhances travel efficiency and minimizes wait times. To our knowledge, this is the first system that combines real-life video feeds, vehicle tracking and a large language model based code generator to create synthetic traffic scenarios and optimize traffic signal plans in SUMO.

1 INTRODUCTION

Simulating traffic scenarios is crucial to understand and optimize urban traffic systems. Tools like **SUMO**, **VISSIM**, and **CityFlow** are commonly used for simulating traffic scenarios (Li et al., 2024) and employing reinforcement learning in traffic signal optimization (Wei et al., 2018), (Oroojlooy et al., 2020) while simulators like **CARLA** and **LGSVL** are being used for autonomous driving research in modern cities (Dosovitskiy et al., 2017), (Rong et al., 2020). Researchers usually use **CARLA** and **LGSVL** to simulate both normal and adverse scenarios (Wei et al., 2024) and then evaluate how autonomous vehicles react under risky driving conditions (Zhang et al., 2024) (Peng et al., 2025). Each of these tools has its own strengths. Among these, **SUMO** (Simulation of Urban MObility) stands out for its microscopic modeling of traffic dynamics, compatibility with real-time data, and open-source accessibility. It is widely used for urban traffic planning, testing signal control strategies, and urban mobility research (Oroojlooy et al., 2020). However, one of the main challenges with such microscopic simulation tools is that one cannot use this software without having any prior knowledge of traffic simulations and traffic-related details. One needs to know how to define road geometry, traffic logic, vehicle behaviors, vehicle types, routes for each vehicle, etc. (Algherbal & Ratrou, 2025).

SUMO allows users to create and simulate complex road networks, define traffic signal logic, test vehicle behaviors, and analyze traffic patterns under diverse edge-case scenarios. Recent studies have shown to integrate large language models into **SUMO** simulators for studying traffic scenarios, congestion and dynamic rerouting. Leveraging large language models **ChatSumo** (Li et al., 2024) and **Sumo-MCP** (Ye et al., 2025) automate traffic simulation through network generation, traffic flow generation.

Congestion in intersections poses a significant challenge in traffic management. Optimization of Traffic Signal Logic (TSC) has been proven to be an effective solution in traffic systems. **SUMO** is an open-source platform for modeling and testing various TSC plans, enabling researchers to observe how optimized signal plans with timing adjustments and real-time interventions can reduce delays and improve traffic flows. In **SUMO**, there are various methods for handling traffic signal control. They can be categorized into three groups: traditional rule-based methods, reinforcement learning (RL) methods, and large language model (LLM)-based approaches. Traditional rule-based methods, such as Fixed-Time Control, Webster’s Method (Ali et al., 2021), Self-Organizing Traffic Light Control (SOTL) (Cools et al., 2008), and Max-Pressure Control (Zoabi & Haddad, 2024), use predefined logic or an algorithm to determine or switch signal phases and their durations. However, they cannot perform well in all the diverse congestion scenarios. RL-based TSC is an approach where the model studies traffic signal control actions by experimenting with the traffic environment,

learns to optimize its control actions after getting feedback dynamically in the form of a reward function (Wei et al., 2019). However, like traditional methods, RL-based methods are not without limitations as well. Although it performs well in dynamic situations, it is unable to perform in unforeseen situations. It cannot handle situations like a broken car in the intersection, a road blockage due to harsh weather, or an emergency vehicle. So Traditional and RL based methods face overfitting issues due to being trained on specific traffic patterns.

To overcome situations like these, traffic signal control researchers nowadays are shifting their focus to LLM-based approaches. LLM-based approaches can handle any complex traffic scenarios that require critical human reasoning, which is again difficult to achieve with traditional rule-based or RL methods. Unlike pre-trained RL models that depend mostly on environment-specific training data, LLMs can generalize across diverse scenarios and respond accordingly after it has reasoned over high-level traffic patterns. Some significant recent works, such as *LLM-Assisted Light* Wang et al. (2024a), *CollMLight* Yuan et al. (2025), have made remarkable use of LLMs in managing traffic congestion by implementing context-aware strategies. Inspired by recent LLM-based works, we study a different but complementary problem: leveraging LLMs for both scenario generation and offline signal plan generation at a single four-way intersection. Specifically, our contributions are: (a) SUMO-compatible scenario generation from natural language prompts, enabling both normal and edge-case conditions; and (b) offline, congestion-aware signal plan generation, where simulation is first run, congestion is assessed, and then the LLM generates static timing plans for a subsequent run.

This framing positions our system as an offline code-generation pipeline, distinct from optimization-oriented approaches. Our approach consists of the following key components:

- **Real-time traffic data collection:** We Collect live traffic video feeds from a four-way urban intersection and apply deep learning-based detection and tracking algorithms to extract vehicle trajectories. These trajectories are aggregated into phase-level occupancy signals.
- **LLM driven Traffic scenario generation:** Unlike prior LLM-TSC papers which only focused on signal control, we use LLMs to create both everyday and edge-case traffic scenarios. Then we simulate these scenarios in SUMO using extracted traffic data.
- **Injecting congestion events into traffic data:** We introduce synthetic congestion events, such as the sudden arrival of a large number of vehicles into normal traffic data extracted from a real-time feed. Then the performance of the LLM to handle the congestion resulted from these events is evaluated.
- **LLM-guided offline traffic signal plans:** LLM is prompted to generate Python code that updates traffic signal logic files based on which phases are maximum in occupancy calculation. It creates different plans for different phase groups. Our LLM generated plans are offline and static which are then evaluated in the next simulation run, rather than optimized online during runtime as in prior TSC work.

2 LITERATURE REVIEW

2.1 SIMULATION IN TRAFFIC RESEARCH

Traffic simulators are becoming increasingly popular in the urban traffic research community. Simulators such as **CARLA** and **LGSVL** (Zhang et al., 2024) are widely used in autonomous driving research for tasks such as crash or collision prevention, safety protocol validation, and sensor robustness testing. On the other hand, microscopic traffic simulators like **SUMO** and **CityFlow** (Yuan et al., 2025) are commonly used to address traffic congestion problems and optimize traffic signal control. Recent research works like **ChatScene** (Zhang et al., 2024) leverage LLMs to generate adversarial, safety-critical scenarios in 3D environments from natural language prompts and help in reducing collisions. It reduces the collision rate to 9%. **LD-Scene** (Peng et al., 2025) uses LLM with Latent Diffusion Models (LDMs) to generate diverse and adversarial driving scenarios from natural language prompts. The frameworks also include Chain-of-Thought code generator and debugger to increase precision and ensure robustness. Other frameworks such as **ChatSim** (Wei et al., 2024) utilize LLMs to run editable, photo-realistic 3D driving scene simulations from natural language prompts. This system focuses on simulating realistic scene videos. **Traffic scene generation from**

natural language description for autonomous vehicles with LLM (Ruan et al., 2024) uses a text-to-traffic scene framework that employs an LLM to generate both standard and critical traffic scenarios for the CARLA simulator based on natural language descriptions. **ChatSUMO** (Li et al., 2024) also uses LLM to generate traffic simulations in SUMO from natural language prompts. The system configures road networks after fetching map data from OpenStreetMap when LLM is given the command to show traffic in any particular city, and then runs SUMO simulations.

2.2 SIGNAL CONTROL IN TRAFFIC RESEARCH

Apart from simulation, TSC is another major focus in urban traffic research. Microscopic simulators are commonly employed to reproduce congestion scenarios and test different solutions under complex environments.

Traditional Rule-Based Method: Webster (Wei et al., 2019) method computes the optimal cycle length and splitting of traffic signal phases at isolated intersections based on traffic flows. Self-Organizing Traffic Light Control (SOTL) (Cools et al., 2008) follows predefined logic to decide whether to switch to another phase or continue the phase based on local vehicle counts. There is a counter that calculates how many vehicles enter the red light zone; if the threshold is crossed, then it decides to switch. Max-Pressure Control (Varaiya, 2013) decides which phase to activate based on the maximum difference between approaching and exiting queue lengths in the intersections to improve the overall traffic travel time. **RL Based Method:** Presslight (Krajzewicz et al., 2002) combines the Reinforcement method with the Max-Pressure theory. The intersection max pressure is key input for the reward function in RL. IntelliLight (Wei et al., 2018) uses the number of vehicles, average waiting time, and lane-wise queue length details. The reward function is mapped to how frequently traffic signals are being changed. UniTSA (Wang et al., 2024b) manages TSC in a Vehicle-to-Everything (V2X) environment. It uses a junction matrix to represent standard intersection and makes use of traffic state augmentation based on intersection symmetry. **LLM Based Method:** LLM Assisted Light (Wang et al., 2024a) optimizes traffic signals dynamically based on maximum occupancies in each phase. For example, when an emergency vehicle suddenly arrives, even if other lanes are slightly congested, that phase will be activated. CoLLMlight (Yuan et al., 2025) optimizes traffic signals dynamically in **Cityflow** simulator. Still, it handles traffic signal control not at a single intersection level, but city-wide.

Inspired by the success and performance of existing works, we have implemented an LLM within the SUMO traffic simulator to handle both normal and abnormal traffic scenarios. Our method is capable of generating routine traffic flows and scenes, as well as simulating congested scenarios. Using the reasoning capabilities of the LLM, the system effectively mitigates congestion issues and manages special traffic situations that often do not fall into the capacities of traditional RL-based methods or even other existing LLM approaches. Unlike other LLM-based methods, our approach does not decide which phase to switch or activate; It rather changes traffic signal plans entirely by generating the simulation code.

3 METHODOLOGY

3.1 OVERVIEW

Our methodology makes use of an LLaMA 3.2 1B model, fine-tuned using QLoRA framework, to generate and optimize traffic scenarios at a four-way intersection in the SUMO simulator. The LLM generates Python code to create both everyday and edge case traffic scenarios and then produces optimized traffic signal plans to reduce congested traffic movements and improve traffic flow.

3.2 INTERSECTION MODELING

To replicate realistic traffic dynamics, we model twelve directional movements corresponding to straight-through and left-turn flows from the four arms of the intersection (North, South, East, West). These movements are categorized into eight *lane groups*:

- **NT, ST, ET, WT** — straight-through movements from north, south, east, and west.
- **N-L, S-L, E-L, W-L** — left-turn lanes in north, south, east, and west.

For traffic scenario generation and signal control, these lane groups are aggregated into four *phase groups*, each corresponding to one controllable signal phase:

- **NS** — North–South straight movements (NT, ST)
- **NS-L** — North–South left turns (N-L, S-L)
- **EW** — East–West straight movements (ET, WT)
- **EW-L** — East–West left turns (E-L, W-L)

3.3 CONGESTION MEASUREMENT

At each simulation timestep t , SUMO reports the instantaneous occupancy of every lane l , denoted $occ_l(t) \in [0, 1]$, where $occ_l(t)$ represents the fraction of lane length occupied by vehicles. We average the occupancies of those lanes.

Lane-group occupancy The occupancy of a lane group g at time t is given by

$$O_g(t) = \frac{1}{|\mathcal{L}_g|} \sum_{l \in \mathcal{L}_g} occ_l(t), \quad (1)$$

where \mathcal{L}_g is the set of lanes in group g .

Time-averaged lane-group occupancy Across a simulation of T timesteps,

$$\bar{O}_g = \frac{1}{T} \sum_{t=1}^T O_g(t). \quad (2)$$

Phase-level occupancy Each phase p consists of multiple lane groups. Its maximum average occupancy is defined as

$$\bar{O}_p = \max_{g \in \mathcal{G}_p} \bar{O}_g, \quad (3)$$

where \mathcal{G}_p is the set of lane groups belonging to phase p . The most congested phase is then selected as

$$p^* = \arg \max_{p \in \mathcal{G}_{phases}} \bar{O}_p. \quad (4)$$

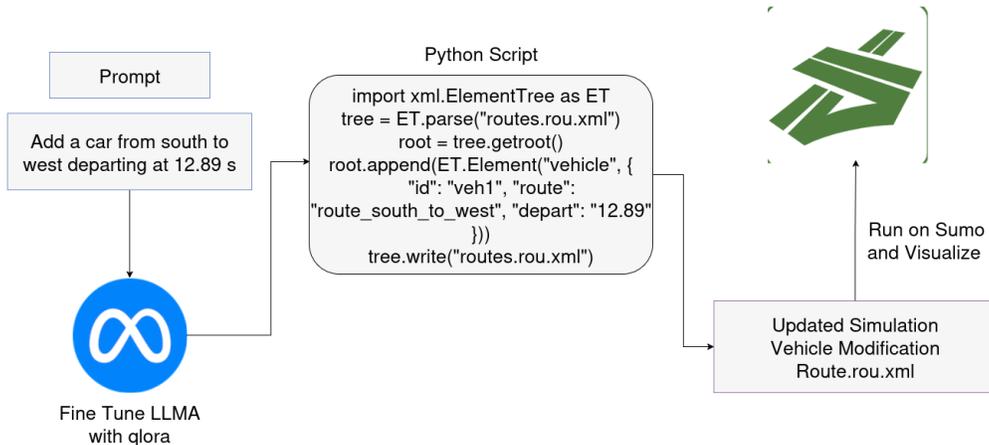


Figure 1: Standard Traffic Scene Generation

3.4 LLM-DRIVEN SCENARIO GENERATION

LLM as a code generator, when a prompt describing a traffic condition is given, it outputs Python scripts that directly write SUMO-compatible XML files. This allows us to create traffic scenarios without manually coding XML. We design three unique scenes:

Scene 1: Normal Traffic Flow The LLM is asked to generate code for normal traffic demand from high-level description (e.g., add a vehicle from north to south) and then generates Python code that automatically updates the route.rou.xml file, specifying unique vehicle ID, departure time, and directional route movement (e.g., north-to-south). Figure 1 illustrates a typical example of such a normal traffic flow, where the LLM-generated script produces routine intersection movements through XML files.

Scene 2: Sudden Vehicle Breakdown In this scenario, the LLM is instructed to place a broken vehicle near the intersection that breaks down while it is taking a left turn. Breakdown can occur in any of the four directions: north, south, east, or west, which can create localized congestion, and we test the robustness of the system to disruptions. For example, given a prompt like A car left at 0.28 s stalls in the west left-turn pocket, backing up traffic, LLM generates python code which updates route.rou.xml file by inserting a special vehicle with a predefined stop position on the west left turn lane at the given time, causing a stall.

Scene 3: Heavy or Medium Traffic Flow High or medium density traffic conditions can be created by giving a command (e.g., Generate a heavy scenario with 231 cars in the network) to LLM to produce python code to generate the route.rou.xml file containing total 231 unique vehicle IDs, randomized departure times, and predefined directional routes, sorted by departure time.

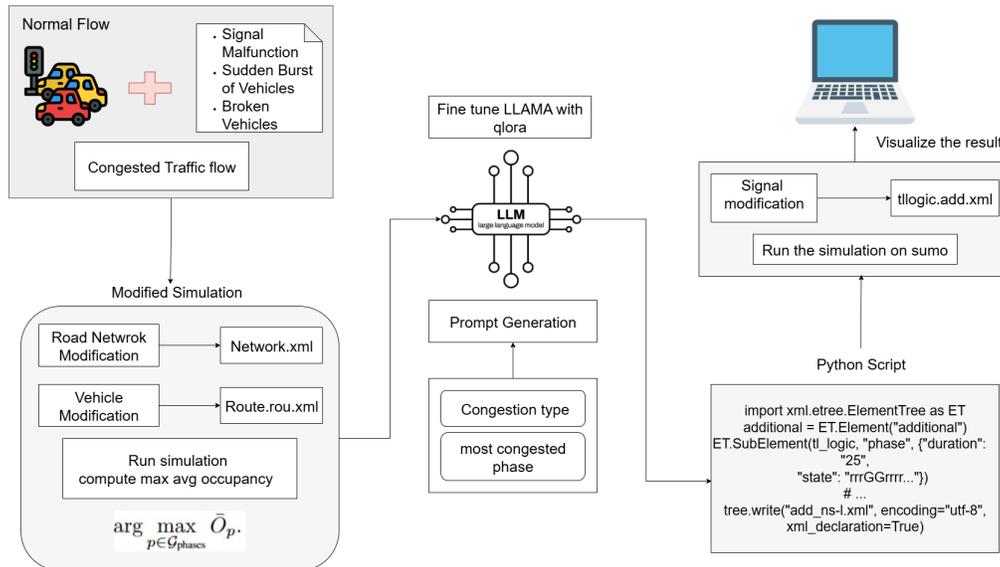


Figure 2: Proposed traffic signal optimization.

3.4.1 SCENE CUSTOMIZATION CATEGORIES

To test the LLM’s performance, we prepared three congestion scenarios. They were scripted manually.

Traffic-Signal Malfunction One or more approaches remain on red light for an abnormally long time, which prevents vehicles from clearing the intersection, thus causing a huge queue of vehicles and congestion.

Sudden Burst of Vehicles A large, short-lived burst of traffic is added along a single direction (e.g., South→North) to test event-driven surges such as before a sports game or during morning office hours. Each of the twelve directions is tested individually, rather than all at once.

3.5 OFFLINE CONGESTION-AWARE SIGNAL PLAN SYNTHESIS

First, we introduce synthetic congestion events—such as signal malfunctions, sudden bursts of vehicles, and stalled vehicles—into otherwise normal traffic flows. These events are encoded by modifying the signal timing in `network.xml` and the vehicle definitions in `route.rou.xml`. The congested scenario is then executed in SUMO, where we compute the average occupancy per phase and identify the phase with the maximum occupancy percentage, along with the corresponding type of congestion. After the simulation, the LLM is prompted with this information (congestion type, most congested phase, and its occupancy percentage) to generate optimized traffic signal plans. The LLM outputs Python code that produces updated, phase-based signal plans in `tllogic.add.xml`. Finally, the scenario is rerun in SUMO to compare traffic performance metrics before and after the application of the LLM-generated plans. The objective of these optimizations is to improve overall traffic flow and reduce the average delay at the intersection. Figure 2 and Algorithm 1 show how LLM handles these three congestion scenarios in our proposed framework.

Green Light Duration Redistribution (Signal Malfunction) To evaluate the LLM’s ability to see if it can recover from severe traffic signal malfunctions, we intentionally designed two extreme failure scenarios—one targeting the entire North-South (NS and NS-L) approach and another one targeting the entire East-West (EW and EW-L) approach. In these cases, the affected approach was intentionally starved by keeping one or more signals red for an abnormally long duration. For Example, the green signal was withheld for 200 seconds (exceeding two full cycles, where each cycle is 90 seconds), followed by only 4 seconds of green time per subsequent cycle in the north approach, which caused both NS and NS-L lanes to suffer. This setup created an extreme level of congestion in the affected lanes. When LLM is commanded, it responds by generating a corrective signal plan that provides longer green durations based on the highest average occupancy percentage observed in each lane group. It produces four timing plans, one for each of the predefined phase groups: NS, NS-L, EW, and EW-L. Each plan adjusts the green time allocation based on congestion severity. Example: Consider a case where the entire East approach is suffering from a traffic signal malfunction. Then it is noticed that the EW phase group suffers the most congestion. An appropriate prompt is passed to the LLM:

Traffic jam in EW because of signal failure, estimated congestion: 48%

In response, the model generates a new Python script. That script outputs a new `tllogic.add.xml` with proper timing plan where green duration for the EW phase is increased most then EW-L phase is given second priority while durations for the other two phase groups are proportionally reduced—allowing the queue to clear faster.

Direction Prioritization During Traffic Surges (Event-driven) In scenarios simulating event-driven surges—such as a burst of vehicles in any of the twelve directional approaches (e.g., South to North, North to East, East to West)—the LLM analyzes congestion by monitoring average phase occupancy levels. Based on the affected phase group (e.g., NS, NS-L, EW, EW-L), it generates a new signal plan with more green time allocated to that phase. For example, if left-turn lanes like N-L (North to East) or S-L (South to West) are congested, an appropriate prompt is passed to the LLM:

NS-L has Highest Congestion with 50% occupancy

Then LLM generates python code to provide more green time to the NS-L phase to help clear the buildup more efficiently. This approach allows the excess queue to clear faster from the overloaded direction, while maintaining a balanced situation in other phases. The LLM always generates four updated signal plans—one for each phase group.

Obstruction-aware Phase Adjustment (Stalled Vehicle) When a vehicle breaks down in a specific lane, such as the left turn lane, the affected lane becomes temporarily inaccessible, as it creates a severe situation by blocking other vehicles behind it from proceeding. This not only delays other vehicles in that lane but can also lead to system-wide congestion due to its effects across connected

lanes. So LLM uses its thinking power to detect the phases that correspond to the affected lanes, then solves the situation by reducing green time in the affected phases and extending green times in unaffected phases until the vehicle is taken away. After the vehicle is taken away, LLM prioritizes the affected phases, extending green time and reducing it for the unaffected phases. LLM’s optimization logic is based on the fact that giving more green time to affected phases is not convenient, as it wastes green time due to vehicles behind the stalled car being unable to proceed. For instance, consider the following scenario where a car breaks down near the east left turn intersection. An appropriate prompt is passed to the LLM:

Generate optimized traffic signal plan for a broken car near the west left-turn intersection.

Model responds by generating a Python script, which in result reduces the green time allocated to that phase (EW-L) and reallocates it to unaffected phases (NS, NS-L, EW) until the vehicle is towed or removed. Once the obstruction is cleared (e.g., the stalled vehicle is towed away), the signal timing is re-adjusted to give priority to the previously blocked phase (EW-L). EW-L phase gets more green time while other phases get less green time.

4 EXPERIMENTS

4.1 MODEL ARCHITECTURE AND FINE-TUNING

4.1.1 LLAMA 3.2 1B WITH QLoRA

We employ supervised fine-tuning (SFT) of the LLaMA 3.2 1B-Instruct model on a dataset of prompt-code pairs (around 5k samples) to generate SUMO-compatible Python code. To reduce memory usage and computational overhead, we are using the QLoRA framework inspired by (Dettmers et al., 2023), which combines **4-bit quantization of the base model** with **LoRA adapters for parameter-efficient fine-tuning**. Specifically, we use the following materials in our fine-tuning pipeline:

4-bit Quantization: When loading the model, we employ a 4-bit compression method (NF4) to downscale the precision of the pretrained model’s weights from 16-bit to 4-bit. Double quantization is enabled, which shrinks the quantization constants even more and ensures lower memory overhead. Pretrained weights are approximately normal $\mathcal{N}(0, \sigma^2)$, we then map them to a fixed distribution by normalizing each block of weights to the range $[-1, 1]$. A codebook with 2^k levels (here $k = 4$) is then derived from the quantiles of the standard normal distribution:

$$q_i = \frac{1}{2} \left(Q_{\mathcal{N}(0,1)}\left(\frac{i}{2^{k+1}}\right) + Q_{\mathcal{N}(0,1)}\left(\frac{i+1}{2^{k+1}}\right) \right), \quad i = 1, \dots, 2^k,$$

where $Q_{\mathcal{N}(0,1)}$ is the quantile function.

LoRA Adapters in QLoRA: Instead of updating the entire 1 B parameters, we add trainable low-rank matrices (rank $r = 8$, scaling $\alpha = 16$) into selected transformer layers: `q_proj`, `k_proj`, `v_proj`, `o_proj`, `gate_proj`, `up_proj`, and `down_proj`. Given a frozen quantized weight matrix W_q , we define ΔW as the low-rank trainable weight update introduced by LoRA:

$$\Delta W = \frac{\alpha}{r} AB,$$

so that the effective weight used during training is

$$W_{\text{eff}} = W_q + \Delta W.$$

Here, A and B are the only trainable parameters, while W_q remains frozen. Thus, QLoRA achieves full-model fine-tuning behavior while only training a few million parameters while keeping the knowledge of the pretrained model intact.

Optimizer: We are using 8-bit paged AdamW optimizer, which minimizes memory usage without damaging the training stability. Gradient updates are applied exclusively to LoRA adapters. Gradient accumulation is used to handle larger batch sizes on limited hardware. **Training Configuration:** The model is trained for 4 epochs with a batch size of 4, gradient accumulation steps of 4, and learning rate of 2×10^{-4} . Gradient checkpointing is enabled to reduce memory usage. The dataset is split into 80% for training and 20% for validation.

4.2 SIMULATION AND EVALUATION

4.2.1 SUMO INTEGRATION

The output from the fine-tuned LLM consists of Python code that outputs XML configuration files in the SUMO traffic simulator. Depending on the type of scenario, the generated code creates one of the following:

route.rou.xml — Defines vehicle routes, departure times, speed, vehicle types, and traffic volumes. This is used for scene generation tasks such as adding vehicles, simulating stalled vehicles, or generating heavy or medium traffic flows.

tllogic.add.xml — Defines custom traffic signal logic. This is created when the LLM is prompted to optimize traffic signal timings in response to congestion scenarios such as signal malfunctions, sudden bursts of traffic, or stalled vehicles.

LLM generates a Python script when it is called. Then this code produces these XML files. Once these XML files are generated, they are loaded into the SUMO to run the simulation to copy real-world traffic conditions. The effectiveness of the LLM-generated control logic is determined by using two well-known traffic performance metrics:

Average Travel Time (ATT): The average time all the vehicles take to travel from their source to their destination.

Average Waiting Time (AWT): The average time all the vehicles wait at intersections.

4.2.2 EVALUATION METRICS UNDER STALLED VEHICLE SCENARIO

To evaluate the outcome of the LLM-generated traffic signal plans, we selected two representative scenarios out of four: one with a broken vehicle near the west intersection and another near the north intersection. Although we have respective scenarios for the east and south approaches, we are focusing on the west and north cases for illustrative purposes. We observed that when the overall vehicle count is low, the differences in traffic metrics before and after applying the LLM-optimized logic are minimal. However, when the traffic flows are very high, lanes remain occupied for some period of time, and the optimized signal plans lead to significant improvements in performance metrics. For both scenarios, the west left turn and the east left turn, we compare ATT and AWT before and after deploying the LLM-generated traffic light adjustments. The results are summarized in Table 1.

Lane	Metric	Before LLM	After LLM
W-L	ATT	932.88	721.27
W-L	AWT	63.59	48.08
N-L	ATT	954.86	640.98
N-L	AWT	63.94	43.05

Table 1: Metrics under stalled vehicle scenario.

4.2.3 EVALUATION METRICS UNDER SUDDEN BURST SCENARIO

To evaluate the effectiveness of the LLM in mitigating congestion caused by sudden bursts of traffic, we conducted high-density traffic simulations from four directions — North to East, North to South, East to West, and East to South — belonging to the NS, NS-L, EW, and EW-L signal phase groups, respectively. Although our system was trained and tested on all 12 directional flows at the intersection, these four directions were selected for detailed analysis. We observed that EW-L routes from East to South and NS-L routes from North to East experienced higher congestion levels and longer delays. These directional approaches are left turn lanes. So vehicles in left turn lanes take more time to clear the intersection than straight lane movements. The delay and wait time in these lanes is always higher than in other lanes.

Phase	Metric	Before LLM	After LLM
NS	ATT	93.97	38.09
NS	AWT	60.89	14.04
NS-L	ATT	1024.42	143.39
NS-L	AWT	81.82	24.47
EW	ATT	96.13	47.55
EW	AWT	61.76	20.82
EW-L	ATT	397.92	41.84
EW-L	AWT	81.53	12.39

Table 2: Comparison of ATT and AWT Across four phases before and after LLM optimization.

4.2.4 EVALUATION METRICS UNDER FAULTY TRAFFIC LIGHT SCENARIO

This extreme traffic signal malfunctioning setup was designed to starve one directional approach for an abnormally longer period, barely giving any green light time to the vehicles to clear the lanes. Thus, a severe level of congestion happens at the intersection. The table below demonstrates how this extreme malfunction led to a dramatic increase in average travel time (ATT) and average waiting time (AWT) before intervention. After applying the LLM-generated signal plans, the system significantly reduced congestion by reallocating green durations based on phase-level occupancy. The difference in ATT and AWT before and after LLM is being used to generate code to optimize traffic flow is much higher, which again reflects the extreme imbalance of the traffic signal configuration in the intersection. Table 3 reports ATT and AWT before and after applying LLM-generated code for signal optimization under extreme congestion conditions.

Phase	Metric	Before LLM	After LLM
NS or NS-L	ATT	2146.73	105.3
NS or NS-L	AWT	75.88	41.8
EW or EW-L	ATT	2261.7	99.7
EW or EW-L	AWT	80.62	38.7

Table 3: Metrics under faulty traffic light scenario.

5 CONCLUSION

Our framework uses a large language model to generate and simulate both normal and adversarial traffic scenarios. Most existing frameworks rely on synthetic data, while a few use real-time data. Some of them use LLMs only to output which phase to activate for traffic signal control. In contrast, our framework derives its dataset from real-world traffic videos, tracks each vehicle along with its direction of movement, speed, and departure time, and incorporates all this information into SUMO simulation, making the simulation more realistic. The LLM then generates code to provide appropriate traffic signal plans to manage congestion. Existing frameworks do not combine all these components into a unified model. From the experimental results, we observe that evaluation metrics like ATT and AWT improve significantly after the LLM is used to manage traffic signal control.

REFERENCES

- Eman A. Algherbal and Nedal T. Ratrout. A comparative analysis of currently used microscopic, macroscopic, and mesoscopic traffic simulation software. *Transportation Research Procedia*, 84:495–503, 2025. ISSN 2352-1465. doi: <https://doi.org/10.1016/j.trpro.2025.03.101>. URL <https://www.sciencedirect.com/science/article/pii/S2352146525001498>. Smart Mobility and Logistics Ecosystems.
- Muzamil Eltejani Mohammed Ali, Akif Durdu, Seyit Alperen Celtek, and Alper Yilmaz. An adaptive method for traffic signal control based on fuzzy logic with webster and modified

- 486 webster formula using sumo traffic simulator. *IEEE Access*, 9:102985–102997, 2021. doi:
487 10.1109/ACCESS.2021.3094270.
- 488
- 489 City of Bellevue, Washington. Bellevue Traffic Video Dataset. GitHub repository, 2017.
490 <https://github.com/City-of-Bellevue/TrafficVideoDataset>.
- 491 Seung-Bae Cools, Carlos Gershenson, and Bart D’Hooghe. *Self-Organizing Traffic Lights:
492 A Realistic Simulation*, pp. 41–50. Springer London, London, 2008. ISBN 978-1-
493 84628-982-8. doi: 10.1007/978-1-84628-982-8_3. URL [https://doi.org/10.1007/
494 978-1-84628-982-8_3](https://doi.org/10.1007/978-1-84628-982-8_3).
- 495
- 496 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. QLoRA: Efficient finetuning
497 of quantized LLMs. In *Thirty-seventh Conference on Neural Information Processing Systems*,
498 2023. URL <https://openreview.net/forum?id=OUIFPHEgJU>.
- 499 Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An
500 open urban driving simulator. In *Conference on robot learning*, pp. 1–16. PMLR, 2017.
- 501
- 502 Daniel Krajzewicz, Georg Hertkorn, Christian Feld, and Peter Wagner. Sumo (simulation of urban
503 mobility); an open-source traffic simulation. pp. 183–187, 01 2002. ISBN 90-77039-09-0.
- 504 Shuyang Li, Talha Azfar, and Ruimin Ke. Chatsumo: Large language model for automating traffic
505 scenario generation in simulation of urban mobility. *IEEE Transactions on Intelligent Vehicles*,
506 pp. 1–12, 2024. doi: 10.1109/TIV.2024.3508471.
- 507
- 508 Afshin Oroojlooy, Mohammadreza Nazari, Davood Hajinezhad, and Jorge Silva. Attend-
509 light: Universal attention-based reinforcement learning model for traffic signal control. In
510 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu-
511 ral Information Processing Systems*, volume 33, pp. 4079–4090. Curran Associates, Inc.,
512 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/
513 file/29e48b79ae6fc68e9b6480b677453586-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/29e48b79ae6fc68e9b6480b677453586-Paper.pdf).
- 514 Mingxing Peng, Yuting Xie, Xusen Guo, Ruoyu Yao, Hai Yang, and Jun Ma. Ld-scene: Llm-guided
515 diffusion for controllable generation of adversarial safety-critical driving scenarios, 05 2025.
- 516 Guodong Rong, Byung Hyun Shin, Hadi Tabatabaee, Qiang Lu, Steve Lemke, Mārtiņš Možeiko,
517 Eric Boise, Geehoon Uhm, Mark Gerow, Shalin Mehta, et al. Lgsvl simulator: A high fidelity
518 simulator for autonomous driving. In *2020 IEEE 23rd International conference on intelligent
519 transportation systems (ITSC)*, pp. 1–6. IEEE, 2020.
- 520
- 521 Bo-Kai Ruan, Hao-Tang Tsui, Yung-Hui Li, and Hong-Han Shuai. Traffic scene generation from
522 natural language description for autonomous vehicles with large language model. *arXiv preprint
523 arXiv:2409.09575*, 2024.
- 524 Pravin Varaiya. Max pressure control of a network of signalized intersections. *Transportation Re-
525 search Part C: Emerging Technologies*, 36:177–195, 2013. ISSN 0968-090X. doi: [https://doi.
526 org/10.1016/j.trc.2013.08.014](https://doi.org/10.1016/j.trc.2013.08.014). URL [https://www.sciencedirect.com/science/
527 article/pii/S0968090X13001782](https://www.sciencedirect.com/science/article/pii/S0968090X13001782).
- 528
- 529 Maonan Wang, Aoyu Pang, Yuheng Kan, Man-On Pun, Chung Shue Chen, and Bo Huang. Llm-
530 assisted light: Leveraging large language model capabilities for human-mimetic traffic signal
531 control in complex urban environments. *arXiv preprint arXiv:2403.08337*, 2024a.
- 532 Maonan Wang, Xi Xiong, Yuheng Kan, Chengcheng Xu, and Man-On Pun. Unitsa: A universal
533 reinforcement learning framework for v2x traffic signal control. *IEEE Transactions on Vehicular
534 Technology*, 73(10):14354–14369, 2024b. doi: 10.1109/TVT.2024.3403879.
- 535
- 536 Hua Wei, Guanjie Zheng, Huaxiu Yao, and Zhenhui Li. Intellilight: A reinforcement learning ap-
537 proach for intelligent traffic light control. In *Proceedings of the 24th ACM SIGKDD international
538 conference on knowledge discovery & data mining*, pp. 2496–2505, 2018.
- 539 Hua Wei, Guanjie Zheng, Vikash Gayah, and Zhenhui Li. A survey on traffic signal control methods.
arXiv preprint arXiv:1904.08117, 2019.

540 Yuxi Wei, Zi Wang, Yifan Lu, Chenxin Xu, Changxing Liu, Hao Zhao, Siheng Chen, and Yanfeng
 541 Wang. Editable scene simulation for autonomous driving via collaborative llm-agents. *2024*
 542 *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 15077–15087,
 543 2024. URL <https://api.semanticscholar.org/CorpusID:267548034>.

544 Chenglong Ye, Gang Xiong, Junyou Shang, Xingyuan Dai, Xiaoyan Gong, and Yisheng Lv. Sumo-
 545 mcp: Leveraging the model context protocol for autonomous traffic simulation and optimization.
 546 *arXiv preprint arXiv:2506.03548*, 2025.

547 Zirui Yuan, Siqi Lai, and Hao Liu. Collmlight: Cooperative large language model agents for
 548 network-wide traffic signal control. *arXiv preprint arXiv:2503.11739*, 2025.

549 Jiawei Zhang, Chejian Xu, and Bo Li. Chatscene: Knowledge-enabled safety-critical scenario gener-
 550 ation for autonomous vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
 551 *and Pattern Recognition (CVPR)*, pp. 15459–15469, June 2024.

552 Razi Zoabi and Jack Haddad. A modified pressure model for max-pressure traffic signal control.
 553 In *2024 European Control Conference (ECC)*, pp. 3370–3377, 2024. doi: 10.23919/ECC64448.
 554 2024.10590907.

555 A DESIGN SPACE COMPARISON

Paper	Data	Task	Propose Signal Plan	Metrics used for Prompt
ChatSUMO	Synthetic	Scenario Generation	No	None
SUMO-MCP	Synthetic	Scenario Generation	No	None
LLM-Assisted Light	Real+Synthetic	Signal Control	Yes	Current phase; Current occupancy
CoLLMLight	Synthetic	Signal Control	Yes	Current phase; Queue length
Ours	Real+Synthetic (video feed)	Scenario and Signal Generation	Yes	Max-occupancy phase; Max avg. occupancy

556 Table 4: Comparison with other LLM-based papers

557 B ADDITIONAL RESULTS

558 B.1 PROMPT-RESPONSE EXAMPLES IN TRAFFIC SCENARIO AND CONTROL GENERATION

559 Normal Traffic Flow

560 Prompt

561 add 3 cars, from west to south, east to north, west to east, departing at 92.79, 95.46, 95.69
 562 seconds

563 Answer

```
564 import xml.etree.ElementTree as ET
565 XML_PATH = r"C:\path\to\routes.rou.xml"
```

```

594
595 def next_ids(root, n):
596     base = max((int(v.get("id").lstrip("veh")) for v in root.findall
597                ("vehicle")),
598                default=-1) + 1
599     return [f"veh{base+i}" for i in range(n)]
600
601 def add_and_sort(routes, departs):
602     tree = ET.parse(XML_PATH)
603     root = tree.getroot()
604     ids = next_ids(root, len(routes))
605     for vid, r, d in zip(ids, routes, departs):
606         root.append(ET.Element("vehicle", {
607             "id": vid,
608             "route": r,
609             "depart": f"{d:.2f}".rstrip("0").rstrip(".")
610         }))
611     vs = sorted(root.findall("vehicle"), key=lambda v: float(v.get("
612     depart")))
613     for v in vs:
614         root.remove(v)
615         root.append(v)
616     tree.write(XML_PATH, encoding="utf-8", xml_declaration=True)
617
618 if __name__ == '__main__':
619     ROUTES = ['route_west_to_south',
620              'route_east_to_north',
621              'route_west_to_east',
622              ]
623     DEPARTS = [92.79, 95.46, 95.69]
624     add_and_sort(ROUTES, DEPARTS)
625
626
627

```

Heavy/Medium Traffic Flow

Prompt

I want to see heavy traffic flow of 259 vehicles.

Answer

```

630
631 import random, xml.etree.ElementTree as ET
632 from pathlib import Path
633 OUT = Path(r"route.rou.xml")
634 N_VEH = 259
635 SIM_END = 1000
636 SEED = 42
637 random.seed(SEED)
638 ROUTES = [('route_north_to_east', 'north_to_center center_to_east'),
639           ('route_north_to_south', 'north_to_center center_to_south'), ('
640           route_north_to_west', 'north_to_center center_to_west'), ('
641           route_east_to_west', 'east_to_center center_to_west'), ('
642           route_east_to_south', 'east_to_center center_to_south'), ('
643           route_east_to_north', 'east_to_center center_to_north'), ('
644           route_south_to_west', 'south_to_center center_to_west'), ('
645           route_south_to_north', 'south_to_center center_to_north'), ('
646           route_south_to_east', 'south_to_center center_to_east'), ('
647           route_west_to_south', 'west_to_center center_to_south'), ('
648           route_west_to_north', 'west_to_center center_to_north'), ('
649           route_west_to_east', 'west_to_center center_to_east')]
650 root = ET.Element("routes")

```

```

648
649 for r_id, edges in ROUTES:
650     ET.SubElement(root, "route", id=r_id, edges=edges)
651     vehicles = []
652     for i in range(N_VEH):
653         depart = round(random.uniform(0, SIM_END), 2)
654         routeID = random.choice(ROUTES)[0]
655         v = ET.Element("vehicle")
656         v.set("id", f"veh{i}")
657         v.set("route", routeID)
658         v.set("depart", f"{depart}")
659         vehicles.append(v)
660     for v in sorted(vehicles, key=lambda e: float(e.get("depart"))):
661         root.append(v)
662     ET.ElementTree(root).write(OUT, encoding="utf-8", xml_declaration=
663         True)
664     print(f" wrote {OUT} with {N_VEH} vehicles (0{SIM_END} s)")

```

Sudden Vehicle Breakdown

Prompt

A car left at 5.48 s stalls in the west left-turn pocket, backing up traffic.

Answer

```

673 import xml.etree.ElementTree as ET, traci
674 from pathlib import Path
675 from sumolib import checkBinary
676 CFG_FILE, ROUTE_FILE = r"sumo_config.sumocfg", Path(r"route.rou.xml"
677 )
678 BROKEN_ROUTE, BROKEN_LANE = "route_west_to_north", "west_to_center_2
679 "
680 DEPART, STALL, END = 5.48, 800, 2300
681 root = ET.parse(ROUTE_FILE).getroot()
682 for n in root.findall('vehicle[@id="broken"]'):
683     root.remove(n)
684 v = ET.SubElement(root, "vehicle",
685     id="broken", route=BROKEN_ROUTE, depart=str(DEPART
686     ),
687     color="255,120,0")
688 ET.SubElement(v, "stop", lane=BROKEN_LANE, pos="82", duration="
689     999999")
690 others = [e for e in root if e.tag!= "vehicle"]
691 vehicles = sorted(root.findall("vehicle"),
692     key=lambda x: float(x.get('depart', '0')))
693 root[:] = others + vehicles
694 ET.ElementTree(root).write(ROUTE_FILE, encoding="utf-8",
695     xml_declaration=True)
696 cmd = [checkBinary("sumo-gui"), "-c", CFG_FILE,
697     "--begin", "0", "--end", str(END),
698     "--collision.action", "none",
699     "--time-to-teleport", "-1"]
700 traci.start(cmd)
701 tow_time = DEPART + STALL
702 while True:
703     traci.simulationStep()
704     t = traci.simulation.getTime()
705     if t >= tow_time and "broken" in traci.vehicle.getIDList():
706         traci.vehicle.remove("broken")

```

```

702
703     if t >= END:
704         break
705     traci.close()
706

```

Traffic Signal Malfunctioning

Prompt

Traffic jam in NS-L because of signal failure, estimated congestion: 48%

Answer

```

716 import xml.etree.ElementTree as ET
717 additional = ET.Element("additional")
718 tl_logic = ET.SubElement(additional, "tlLogic", {
719     "id": "center",
720     "type": "static",
721     "programID": "sudden_ns-1",
722     "offset": "0"
723 })
724 ET.SubElement(tl_logic, "phase", {"duration": "25", "state": "
725     rrrGGrrrrrrrrGGrrrrr"})
726 ET.SubElement(tl_logic, "phase", {"duration": "5", "state": "
727     rrryyrrrrrrrryyrrrrr"})
728 ET.SubElement(tl_logic, "phase", {"duration": "30", "state": "
729     GGGGGrrrrrGGGGGrrrrr"})
730 ET.SubElement(tl_logic, "phase", {"duration": "5", "state": "
731     yyyyyrrrrryyyyyrrrrr"})
732 ET.SubElement(tl_logic, "phase", {"duration": "20", "state": "
733     rrrrrGGGGrrrrrGGGGG"})
734 ET.SubElement(tl_logic, "phase", {"duration": "5", "state": "
735     rrrrryyyyyrrrrryyyy"})
736 tree = ET.ElementTree(additional)
737 tree.write("add_ns-1.xml", encoding="utf-8", xml_declaration=True)
738

```

Sudden Burst of Vehicles

Prompt

EW has Highest Congestion with 50% occupancy

Answer

```

745 import xml.etree.ElementTree as ET
746 root = ET.Element("additional")
747 tl = ET.SubElement(root, "tlLogic", id="center", type="static",
748     programID="solution_faulty_lights_ew", offset="0")
749 ET.SubElement(tl, "phase", duration="45", state="
750     rrrrrGGGGrrrrrGGGGG")
751 ET.SubElement(tl, "phase", duration="10", state="
752     rrrrryyyyyrrrrryyyy")
753 ET.SubElement(tl, "phase", duration="25", state="
754     GGGGGrrrrrGGGGGrrrrr")
755 ET.SubElement(tl, "phase", duration="10", state="
756     yyyyyrrrrryyyyrrrrr")

```

```

756
757 tree = ET.ElementTree(root)
758 tree.write("tlLogic_output.add.xml", encoding="utf-8",
759           xml_declaration=True)
760

```

Stalled Vehicle near left lane in intersection

Prompt

Generate optimized traffic signal plan for a broken car near the east left-turn intersection.
Scenario 10

Answer

```

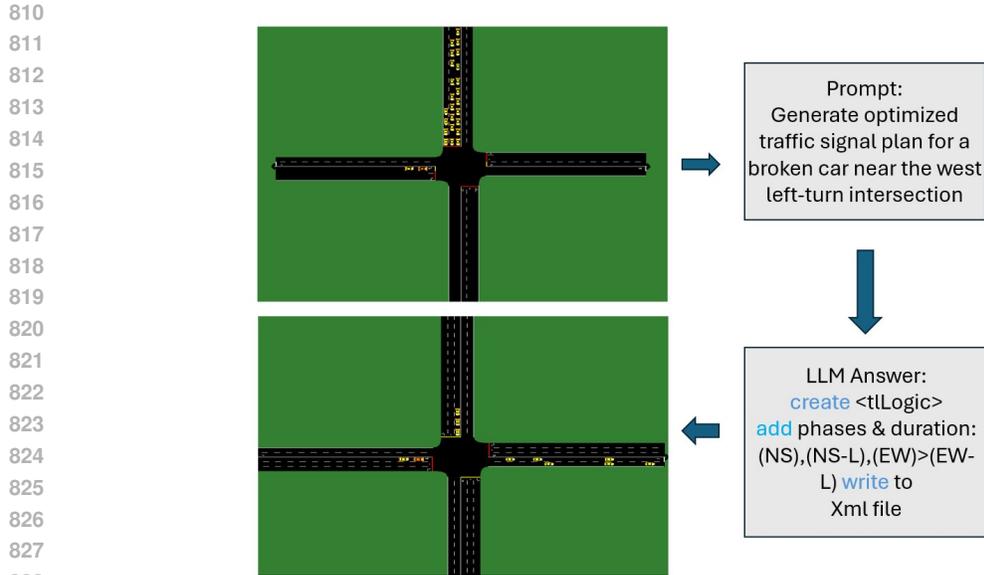
770
771 import xml.etree.ElementTree as ET
772 root = ET.Element("additional")
773 broken_logic = ET.SubElement(root, "tlLogic", id="center", type="
774 static", programID="broken_eleft", offset="0")
775 for duration, state in [
776     ("25", "rrrrrGGGrrrrrrGGGGG"), ("5", "rrrrryyyrrrrrryyyyy"), (
777     "1", "rrrrrGGGGrrrrrrGGGGG"), ("4", "rrrrryyyrrrrrryyyyy"), ("
778     50", "GGGGrrrrrrGGGGrrrrr"), ("5", "yyyyyrrrrrryyyyyrrrrr")
779 ]:
780     ET.SubElement(broken_logic, "phase", duration=duration.strip('')
781     ), state=state.strip(''))
782 faulty_logic = ET.SubElement(root, "tlLogic", id="center", type="
783 static", programID="faulty_EL", offset="0")
784 for duration, state in [
785     ("50", "rrrrrGGGGrrrrrrGGGGG"), ("10", "rrrrryyyrrrrrryyyyy"),
786     ("20", "GGGGrrrrrrGGGGrrrrr"), ("10", "yyyyyrrrrrryyyyyrrrrr")
787 ]:
788     ET.SubElement(faulty_logic, "phase", duration=duration.strip('')
789     ), state=state.strip(''))
790 tree = ET.ElementTree(root)
791 tree.write("generated_broken_eleft.add.xml", encoding="utf-8",
792           xml_declaration=True)
793

```

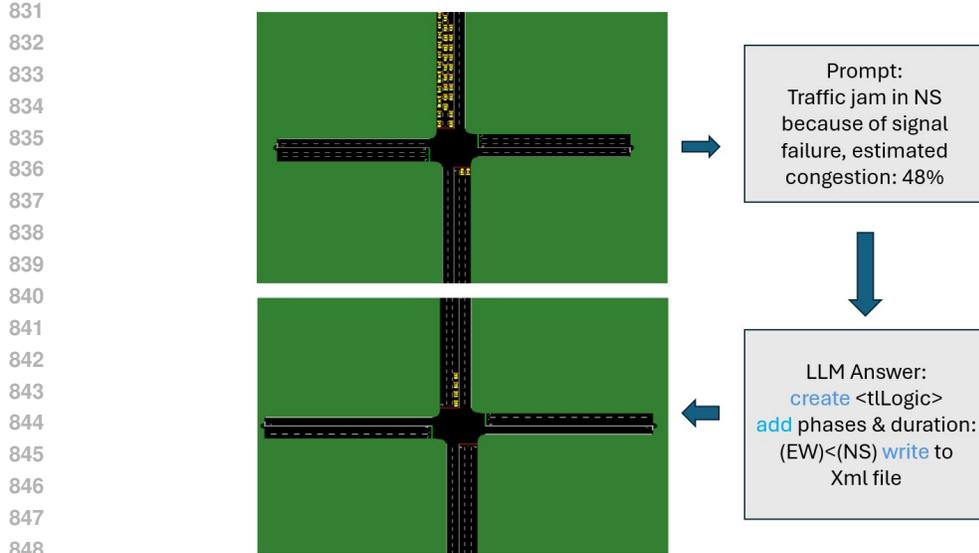
B.2 VISUALIZATION OF THE RESULTS

Stalled Vehicle Scenario At $t = 357$ s, the stalled vehicle (orange car) in the EW-L lane is still present, so it cannot discharge. The LLM-generated Python script responds by writing an updated `tlLogic.add.xml` that temporarily reduces green allocated to EW-L and reallocates it to NS, NS-L, and EW to keep the intersection flowing. This behavior prioritizes phases that can actually serve demand while the obstruction persists. Once the vehicle is removed, the timing plan rebalances to restore priority to EW-L, granting it additional green while proportionally shortening the others. We therefore compare outcomes at $t = 357$ s (vehicle not yet towed) with and without the LLM plan in figure 3

Traffic Signal Malfunction At $t = 577$ s, the northbound approach is heavily congested without the LLM plan, whereas with the LLM plan at the same timestamp the northbound vehicle count drops markedly in figure 4. Given these congestion inputs, the LLM generates Python that writes `tlLogic.add.xml` to extend green for NS (then NS-L) while proportionally shortening the other phases, yielding the observed reduction.



829 Figure 3: comparison without LLM vs. with LLM-generated timing plan in broken car scenario.



850 Figure 4: Traffic signal malfunctioning: comparison without LLM vs. with LLM-generated timing plan

853 C DATASET AND ALGORITHM

855 C.1 DATASET COLLECTION

856
857 We utilized a real-world traffic dataset collected from the Bellevue intersection in Washington. City of Bellevue, Washington (2017). Vehicle trajectories were extracted using YOLO object detection and the DeepSORT tracking algorithm, which allowed us to determine the number of vehicles entering from each approach, their respective movements, average speed, and departure times. Subsequently, these trajectories were integrated into the SUMO traffic simulator to follow realistic traffic flow conditions.

860
861
862
863 The dataset spans eight days and was segmented into six temporal intervals: 1-minute, 2-minute, 3-minute, 4-minute, 5-minute, and 6-minute durations. For sudden bursts of vehicles, we used

864 1-minute and 2-minute duration videos. For faulty traffic signal scenarios, we used 2-minute, 3-
 865 minute, 4-minute, 5-minute, and 6-minute videos. We trained on 6 days of data, then tested on 2
 866 days of unseen data. Based on this, we introduced a variety of congestion-inducing scenarios to
 867 simulate more complex and realistic urban conditions. These included:

- 868 1. Placement of stalled vehicles in the left lane from north/south/east/west directions to simu-
 869 late road obstructions,
- 870 2. Extreme malfunctioning of the traffic signal, which results in an unusually longer red phase
 871 in one direction, and
- 872 3. Sudden bursts of vehicles during peak hours to emulate event-driven surges.
 873

874
 875 Following the integration of these scenarios, we measured the average occupancy percentage for
 876 each signal phase—defined as the proportion of the lane occupied by vehicles. The average occu-
 877 pancy percentage served as key input for the LLM to assess congestion levels and generate optimized
 878 traffic signal control plans.

879 C.2 ALGORITHM OF THE SYSTEM

880 The workflow monitors traffic phases, computes occupancies, and detects congestion types (stalled
 881 vehicle, malfunction, or burst). It then queries the LLM with occupancy data and congestion type to
 882 generate updated signal plans. Finally, SUMO is rerun with the new plans and performance metrics
 883 are compared.
 884

885 **Algorithm 1** LLM-Guided Traffic Signal Optimization Workflow

```

886 1: Initialize SUMO simulation
887 2: while simulation not finished do
888 3:   Map each lane to its phase group {NS, NS-L, EW, EW-L}
889 4:   Compute per-phase average occupancy  $\bar{O}_p$ 
890 5:    $p^* \leftarrow \arg \max_p \bar{O}_p$ 
891 6:   if stalled vehicle detected then
892 7:     Mark congestion type  $\leftarrow$  STALLED VEHICLE
893 8:     Record towing time, vehicle departure time, affected lane
894 9:   else if signal malfunction detected then
895 10:    Mark congestion type  $\leftarrow$  SIGNAL MALFUNCTION
896 11:  else
897 12:    Mark congestion type  $\leftarrow$  SUDDEN BURST
898 13:  end if
899 14: end while
900 15: Call LLM with ( $\{\bar{O}_p\}$ , congestion type)
901 16: Receive Python code that generates new signal plans
902 17: Run SUMO with new signal plans
903 18: Record and compare metrics (ATT, AWT) before and after LLM
904
```

905
 906 Algorithm 1 illustrates how the system handles each congestion event.
 907
 908
 909
 910
 911
 912
 913
 914
 915
 916
 917