# Decision Support for Traffic Management using Reinforcement Learning

*Author:*
A. Heijnen

*Supervisors:*
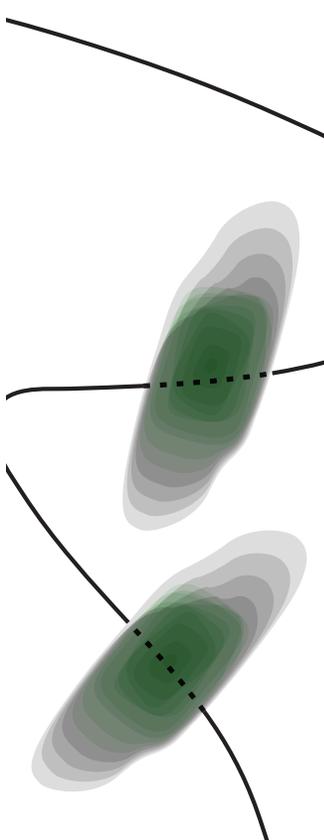M.R.K Mes
W.J.A. van Heeswijk
P.W.M. van Bakel

*Case owners:*
Technolution

*Date:*
February 12, 2025

# Management Summary

This research investigates the application of Deep Reinforcement Learning (DRL) as a decision-support tool with the objective to reduce urban traffic congestion across multiple intersections. As the currently used solution method in traffic management, rule-based decision-making, is time-consuming to implement, lacks scaling, and is difficult to design well, an alternative is required. DRL is chosen as the solution method due to the complex nature of traffic and the long-term effects of actions. The strengths of neural networks and Reinforcement Learning (RL) are combined to learn the long-term effects of actions and find the best solution given a traffic state. The main research question is "**How suitable is deep reinforcement learning for providing real-time decision support to traffic engineers for reducing congestion across multiple intersections in an urban traffic network?**"

As the aim is to provide decision support to traffic engineers, the outcome of the DRL algorithm has two requirements: interpretability and infrequent adjustments. This research adopts a macro-level approach that adjusts traffic flows over 15-minute intervals, unlike micro-level optimizations that alter traffic light phases every few seconds. The selected method, perimeter control, regulates traffic that enters and exits a defined city center region by modifying traffic light phase durations at its perimeter, being the outside edges of the region. This strategy reduces congestion in the city center, maintaining fluent traffic flow with reduced average travel times.

However, real-world traffic is rarely uniform, resulting in potential inefficiencies such as wasted green time at certain traffic lights. To address this, a secondary layer of control is introduced, max pressure. This micro-level technique balances traffic queues at individual intersections by considering incoming traffic flow and the available space at destinations. Together, these two layers form an integrated control system: Perimeter control for strategic adjustments combined with max pressure for micro-level control for localized optimization.

First, this study simulates a 25-intersection grid network using the SUMO traffic simulator, representing a real city scenario. Supervised learning is then used to train a neural network capable of predicting average traffic speed over 7.5-minute intervals. The results show a prediction accuracy with a mean absolute percentage error (MAPE) of 20–25%.

Second, iterative supervised learning is used to choose the best action given a traffic state. The newly generated data is then used to fine-tune the neural network with the newfound results. After five iterations of using the trained neural network to choose the traffic light policy, the average speed is increased by 5-9% with respect to the initial performance of static random policies. When compared to a better standard static policy the mean speed is increased by only 0-3%. Although the mean speed is increased by this approach, the results are bad since the policy it is compared to involves no intelligence and should thus be significantly outperformed. On the positive side, there is still potential for further improvements as more iterations are done. The process of iterative supervised learning is also computationally intense making it unrealistic to perform a large additional amount of iterations. Fine-tuning the original neural network does improves the prediction accuracy, with a newfound MAPE of 8–10%.

Lastly, DRL is used to find the best action given a traffic state. This is first done without max pressure and thus no local traffic light optimization. All traffic lights are changed with the same amount according to the action found by the DRL algorithm. As more iterations are done, learning can be observed, as the reward increases. Despite this, the increase in performance plateaus slightly below the mean speed of the standard traffic light settings. This means that the best-found solution of the DRL algorithm performs worse than making no changes, which is also a possible action for the DRL algorithm to select. This leads to the belief that the DRL algorithm is not able to fully understand the impact of an action on the environment or the long-term effects. Supervised learning shows that predicting future traffic is achievable. Nonetheless, choosing actions to improve the average speeds leads to relatively small improvements.

Reflecting on the research question: In theory, DRL should be able to provide decision support to traffic engineers and in consequently reduce congestion in the city. However, the implementation of DRL is challenging and unsuccessful during this research. Alternatively, supervised learning using neural networks is found to be effective in predicting traffic. This could be used as a form of decision support for traffic engineers.

Key findings and recommendations:

- Supervised learning using neural networks is a promising tool for predicting complex data such as traffic. With additional fine-tuning of the neural network a MAPE of 8-10% is achieved. Supervised learning can also be used to make actions increasing the average speed by 0-3% compared to a standard static policy.

- Deep reinforcement learning (DRL) should in theory be a useful tool for traffic management but implementation is difficult due to the black-box nature of neural networks. In this study, DRL is not able to outperform the base traffic light policy. If DRL is the desired solution method for a problem, sufficient time and expertise is required for successful application.

- For gathering data, both supervised learning and DRL are computationally intense. Once both techniques are trained, the time required to find a good solution is short. Hence, both techniques are still applicable in cases that need fast response times.

- The most useful artificial intelligence tool identified during this research is neural networks, which are used to predict future traffic. This can assist in providing decision support to traffic engineers and could be a promising service for Technolution to provide.

- To investigate techniques to reduce computational time future research is recommended by the author. As well as exploring alternative reward structures, state formulations, and actions to improve the DRL approach.

## Acknowledgements

I am grateful for the opportunity that Technolution has presented with this research. I always enjoy coming to the office and talking to helpful, friendly colleagues. Completing this research would not have been possible without the help of my supervisor from Technolution, Tijs van Bakel. I offer my sincere appreciation for all the help and useful feedback.

A special thanks to my first- and second supervisor, Martijn Mes and Wouter van Heeswijk, for giving guidance, and thorough feedback during my research. Your input has significantly helped me during this project.

Lastly, I would like to thank my girlfriend, family, and friends for the review, encouragement, and support that helped me in completing this research.

# Contents

# List of Figures

# Glossary

| | |
|---|---|
| Ant Colony Optimization | A bio-inspired algorithm that mimics the foraging behavior of ants. Ants leave pheromones on paths while searching for food, and the algorithm simulates this by favoring paths with stronger pheromone concentrations. Over time, the algorithm converges to optimal or near-optimal solutions. |
| Deep Loosely Coupled Q-network (DLCQN) | Type reinforcement learning which allows agents to act independently or cooperatively depending on circumstances to prevent policy bias. |
| Deep Reinforcement Learning (DRL) | Technique of training software to make decisions using a neural network to get the optimal result. |
| Deep Repeated Update Q-network (DRUQN) | Type of reinforcement learning which updates action values depending on the probability of selecting an action to avoid policy bias. |
| Markov Decision Process (MDP) | Mathematical framework for modeling decision making where outcomes are partially random and controlled by the decision maker. |
| Reinforcement Learning (RL) | Technique of training software to make decisions to get the optimal result. |

**Note:** In this thesis, the terms "Reinforcement Learning (RL)" and "Deep Reinforcement Learning (DRL)" will be used interchangeably, as DRL is just RL using neural networks.

# 1 Introduction

This chapter gives background information about the company, a description of the problem, as well as an approach to solving this problem with the corresponding research questions and subquestions. To conclude, an overview of the research design will be provided.

## 1.1 Company description

Technolution is a Dutch technology company founded in 1987. Technolution specializes in developing innovative, high-tech solutions across various industries, including mobility, energy, and security. The mobility sector focuses on aiding traffic engineers, finding solutions to traffic problems and automatization of processes. Their software MobiMaestro is the most used platform for traffic management in The Netherlands. MobiMaestro provides users with insight, overview, and control regarding environment, traffic, and city management (Technolution Move, 2024).

## 1.2 Research plan

The following subsection will cover the problem description as well as the objective of this research based on the assignment provided by Technolution.

### 1.2.1 Problem description

Technolution would like to have a new system where the users (traffic engineers) can implement their conditions by assigning priorities to the road network. The new system will provide decision support to traffic engineers, making it easier for them to configure cities based on their prioritized objectives. This decision support system would tell the traffic engineer which changes to make to the traffic control system to get the preferred result of the traffic engineer, making it a prescriptive decision support system. Nonetheless, determining the optimal traffic control system settings for all systems presents a complex algorithmic challenge. In large cities, there are hundreds or even thousands of traffic lights. The different traffic situations that are encountered are also very high, making the mathematical model and thus algorithm complex to define and model. As it is not feasible to mathematically compute an optimal solution with so many traffic scenarios, a degree of uncertainty, and long-term effects, an alternative approach is required that can address these challenges without relying on exhaustive calculations of every possible scenario (Inoue et al., 2024).

### 1.2.2 Objective of the research

As Technolution wants to explore the possibility of giving decision support to traffic engineers, a solution method must be interpretable and require infrequent adjustments. In addition, the solution method must address the challenges of managing long-term effects and navigating an extremely large number of scenarios that are computationally infeasible to solve through traditional mathematical iteration.

The primary objective of this research is therefore to design a way to give decision support to traffic engineers that fulfills the aforementioned requirements. A solution method that Technolution is interested in is Reinforcement Learning (RL) for its ability to make sequential decisions under

uncertainty while accounting for long-term effects. By combining RL with deep learning as Deep Reinforcement Learning (DRL), it can use a neural network to generalize and efficiently predict the long-term effects given a state (Li, 2018). This research aims to use the strengths of DRL to provide traffic engineers with interpretable recommendations for managing traffic efficiently.

## 1.3 Research questions and approach

This study investigates the ability to use the theoretical advantages of DRL for traffic optimization while overcoming challenges such as interpretability and scalability. In accordance with the objective of the research, the following main research question will be answered:

> *How suitable is deep reinforcement learning for providing real-time decision support to traffic engineers for reducing congestion across multiple intersections in an urban traffic network?*

To answer the main research question, additional questions and subquestions are formulated. These are as follows:

- How is traffic management used to improve the desired flow of traffic?

  - What are commonly used techniques in city traffic management?
  - What are the challenges in traffic management?

- How has artificial intelligence been applied in traffic management according to the literature?

  - How is supervised learning used for traffic management?
  - What makes RL a suitable solution for city traffic management?
  - What are the main challenges of RL in city traffic management?

First, a literature study will be conducted to gain a better understanding of traffic management and to study the application of artificial intelligence in traffic management. This literature review covers the subquestions of questions one and two. The literature study ends with a short conclusion to answer questions one and two: "How is traffic management used to improve the desired flow of traffic?" & "How has artificial intelligence been applied in traffic management according to the literature?".

- How can the solution method be designed to reduce traffic congestion while being able to provide decision support to traffic engineers?

  - What are the key elements of the traffic congestion problem?
  - How can the problem be formulated as a MDP?
  - How can the traffic congestion problem be solved?
  - How will this solution be given as advice to traffic engineers?
  - What RL method is most suitable for solving this city traffic management problem?

Second, the findings of the literature will be used to model the problem and find a fitting solution. This chapter will include a model of the problem in terms of a Markov decision process (MDP), as well as a solution method that can be applied to solve the problem and the tools needed for this. The chapter on the modeling approach concludes with an answer to question three "How can the solution method be designed to reduce traffic congestion while being able to provide decision support to traffic engineers?".

- How should the RL algorithm be applied to optimize traffic flow?
    - What modeling choices and assumptions need to be made?
    - How will the RL algorithm be verified?
    - How can computational time be kept at a minimum?

Third, after modeling the problem and the solution method, the solution method will be applied. The traffic simulator Simulation of Urban MObility (SUMO) will be used to first simulate traffic and second to understand how traffic policies influence flow under various circumstances. Using the traffic simulator, many phase durations and traffic intensities are simulated to gather data to train a neural network. This trained neural network is then used for iterative supervised learning where the neural network decides the best policy given a traffic state. The newly generated data is also used to fine-tune the existing neural network. After this approach, DRL will be applied to find the best action given a traffic state. The chapter on applying the reinforcement learning algorithm concludes by answering question four: "How should the RL algorithm be applied to optimize traffic flow?".

- How well do supervised learning and the RL algorithm perform?
    - How well can a neural network predict future traffic congestion?
    - How well can a neural network be applied to modify traffic light configurations to improve the mean speed?
    - How does the performance of the RL algorithm compare to the original traffic performance in terms of mean speed?
    - What are the limiting factors for implementation?

Fourth, the performance of the neural network as well as the RL algorithm are analyzed. To do this, the prediction accuracy of the neural network is determined. In addition to evaluating the performance, in terms of mean speed, when the neural network is applied to change the traffic light configurations. The performance of the RL algorithm will be compared to a benchmark policy for validation. These results will answer the subquestions of question five. The chapter on experimental results provides an answer to question five: "How well do supervised learning and the RL algorithm perform?".

Lastly, the conclusion presents the answer to the main research question: "How suitable is deep reinforcement learning for providing real-time decision support to traffic engineers for reducing congestion across multiple intersections in an urban traffic network?", as a final result. This will be done by reflecting on the results of the previous chapters and their research questions and subquestions.

9

## 1.4 Outline of the report

Chapter 2, Literature study, gives insights into the challenges and solution methods within traffic management. As well as how supervised learning and DRL are applied in traffic management. The chapter will conclude by answering the first two research questions. Chapter 3, Approach for modeling the problem and the solution method, includes a model of the problem as well as an in-depth explanation of the approach used to solve this problem. The chapter concludes with an answer to the third research question. Chapter 4, Experimental settings, covers how the solution method is implemented. Including all the steps taken to get towards the final solution. The chapter concludes with an answer to research question four. Chapter 5, Experimental results, covers the results found during this research. This covers the training of the neural network, the iterative supervised learning approach as well as the DRL approach. The chapter will conclude with an answer to the final research question. Chapter 6, Discussion, will discuss the findings of the results including limitations and challenges. Chapter 7, Conclusion, will summarize the research and the found results, rounding off by answering the main research question. The research finishes off with Chapter 8, Recommendations and future research, in which the author outlines potential directions for future research as well as recommendations for Technolution.

## 2 Literature study

This chapter aims to gain a better understanding of the challenges within traffic management and to gain familiarity with solution methods for traffic management problems. Additionally, the chapter is used to gain insight into how supervised learning and reinforcement learning are used in traffic management. The chapter rounds off by answering subquestion 1 and 2: "How is traffic management used to improve the desired flow of traffic?" & "How has artificial intelligence been applied in traffic management according to the literature?".

### 2.1 Challenges and solutions for traffic management

The first step towards solving a traffic management problem is to understand the challenges encountered in traffic management. The major challenge in traffic management is that the number of vehicles increases faster than the infrastructure in place supporting it. This leads to road congestion which has various consequences: Poor road user experience, loss of time, a decrease in productivity, and an increase in pollution (Ouallane et al., 2022). Ouallane et al mention two main principles to combat traffic congestion: Traffic guidance and traffic light control. Traffic guidance includes different methods for finding alternative routes for vehicles with the aim of preventing congestion. Traffic light control includes several methods to improve traffic flow at intersections and thus also aims at the reduction of congestion. Further, Ouallane et al show that researchers tend to propose traffic light control solutions for isolated intersections. However, for the solutions to be efficient, coordination between intersections is important to have an understanding of local fluent flow but also a global fluent flow. Ouallane et al conclude their work by recommending the combination of modern sensor technology to gather data and artificial intelligence to optimize efficient road traffic management.

In the Dutch booklet Handboek verkeersmanagement – Module Regelaanpak (CROW, n.d.) they mention four possible actions a traffic manager can take to solve congestion. These actions are: Redirecting traffic, reducing inflow, promoting outflow, and informing road users. All four options can cause follow-up problems as traffic is changed. Two possible follow-up problems can be overloaded alternative routes and creating new bottlenecks. Consequently, this challenge will have to be addressed as well. Some solutions for redirecting traffic may be a good solution but are not allowed in The Netherlands (CROW, n.d.). An example is that traffic is not allowed to be redirected in a way that disturbs the cycling traffic of kids going to school in urban areas. An additional challenge that traffic engineers have to deal with is that although road users can be informed of upcoming traffic situations, they might not see the sign or simply choose to ignore it. The combination of the above-named challenges suggests the complexity of solving traffic problems and that solutions can require multiple steps.

Ravish and Swamy (2021) did research into intelligence-based algorithms, and computational methods that mimic human or natural intelligence, using artificial intelligence, to solve tasks. Ravish and Swamy found that intelligence-based algorithms, used to mitigate congestion, require a large time to compute. The computational time for Ant Colony Optimization for a network of 142 links is found to be almost 14 hours (Cong et al., 2013). Ravish and Swamy found that genetic algorithms solve some issues such as tuning parameters and handling of objectives and constraints. However, this leads to even higher computational costs than ant colony optimization. Swarm

optimization is another intelligence-based algorithm that is used for traffic optimization. Shi et al. (2021) applied swarm optimization to optimize a single large intersection. They were able to reduce the waiting time by 20%. However, the computational times are again long, ranging from 5 to 14 hours depending on the traffic intensity that is simulated.

A relatively simple technique for optimizing local traffic lights is the use of the max pressure (MP) algorithm (Varaiya, 2013). The algorithm is used for deciding the duration of green time per phase. A phase is defined as a possible combination of green and red lights at an intersection. A commonly used approach to determine phases for traffic lights is the ring-barrier approach. This approach gives eight possible phases as can be seen in Figure 1. These phases can transition into each other. The implementation of the ring barrier approach works as follows: The traffic light cycle starts with phases 1 and 5 concurrently. When phase 1 ends, phase 2 can start, while phase 6 waits for phase 5 to finish. The dotted lines shown in Figure 1 indicate permissive phases. In contrast to protected phases which always have priority, permissive phases are allowed to go if there is a gap but must give priority to conflicting phases. Returning to the working of max pressure, the choice of the green duration of a lane is based on the queue of that lane, which is interpreted as pressure, at the intersection as well as the queues of the links around the intersection (Tsitsokas et al., 2021). This is shown in Figure 2. The lane which has the most pressure is then given green.

Varaiya (2013) mentions three advantages of MP. First, since MP is for local intersection control it only requires the queue length of the lanes at the intersection and the queue lengths of the links. As a result, the communication requirements between intersections are low. There is also no need to change other intersections as a direct result of the action at one intersection. Second, MP adjusts phase durations to a shift in demand without the need to predict or estimate demands. Third, normal MP requires knowledge of turn rations and saturation flow rates. However, by implementing the adaptive MP method, those parameters can be estimated, further simplifying the implementation of max pressure.

For a larger scale of traffic management, perimeter control can be implemented (Tsitsoka et al., 2023). This method "guards" an area enclosed by a perimeter. The guarding behavior can be explained as a perimeter line that moderates the incoming traffic such that a desired traffic outcome in the area is achieved. This can be a limited amount of congestion or a minimum speed. An example of this application is in the Maas tunnel in Rotterdam. Here multiple layers of perimeter control are implemented to prevent the average tunnel speed from dropping below 15 km/h (Technolution Move, 2023). This is achieved by moderating inflow at various intersections surrounding the tunnel. Here the goal is not to improve overall traffic flow but rather to fulfill the must-have safety requirements for driving in the tunnel such that cars are not slower than the smoke or toxic fumes in the tunnel in case of a calamity. The reason for this is that because the tunnel is not allowed to have visible ventilation at the entrance as it is a national monument, it uses longitudinal ventilation. This will push the smoke at a speed of 15 km/h. If traffic were to move more slowly than this, cars could become stuck in the smoke, creating a safety hazard. The PC approach used for the Maastunnel can also be applied to cities. By reducing the in-flowing traffic on the outskirts of the city, the congestion inside the city can be reduced.

Figure 1: 8 phases in the ring barrier traffic signal control (`https://ops.fhwa.dot.gov/publicati ons/fhwahop08024/chapter4.htm`).

## 2.2 Supervised learning in traffic management

When trying to solve a problem that changes over time it is important to know how the problem will evolve. This is also the case for traffic management as Tampubolon and Hsiung (2018) state: Traffic management systems rely on the ability to predict traffic flows. Tampubolon and Hsiung trained a neural network to predict traffic for a day in the future. They use the past 3 days of traffic as well as the historical traffic of that day. They found good performance on training data but the neural network struggled to generalize and would instead overfit on training data. The result is that the validation performance was rather poor. They predict that the small size of the training data is the cause of this. To prevent the neural network from overfitting they added dropout regularization. This technique will randomly choose to ignore a percentage of the neurons. As a neural network trains the neurons can specialize, developing specific roles in the neural network. This can cause neighboring neurons to depend too heavily on this specialization. The result is that complex co-adaptations are formed which causes the model to be overly specific for the training data. By randomly removing neurons this effect can be prevented and instead, neurons will be more independent preventing overfitting (Sanagapati, 2019). After applying the dropout technique the neural network was able to generalize better. Tampubolon and Hsiung found that the availability of data was a great limitation and more data should improve the performance. They also state that the long time required for training a deep neural network on a large amount of data is time-consuming.

A review of the applications of deep learning in traffic management by Patil (2022) mentions that deep learning applications are seeing an increase in uses in traffic management. The applied areas are traffic prediction, object detection and recognition, autonomous vehicles, traffic flow optimization, and parking management. Traffic prediction is the most used application, enabling traffic engineers to anticipate congestion. Despite positive results, Patil states that traffic engineers can struggle to interpret why deep learning models are making certain predictions. In addition, traffic managers experience difficulties adjusting the parameters of the deep learning models to

Figure 2: Illustration of max pressure control.
Illustration of max pressure control. In Case A, the green signal is set in the north-to-south direction. In Case B, the green signal is set in the east-to-west direction (Wei et al. 2019a).

increase performance. This can limit the effective use of deep learning applications. In accordance with results found by Tampubolon and Hsiung, Patil also mentions that the need for a large amount of data is a limitation of deep learning. A solution given by Patil is to continuously train the model on new data to increase accuracy.

## 2.3 RL in traffic management

After being able to predict traffic using a neural network the next step is to perform actions taking into account the long-term effects that are predicted, this is the field of Reinforcement Learning (RL). RL is useful in cases where there is a stochastic environment with sequential decision-making. This is very fitting to traffic as traffic is rapidly changing and decisions will have a long-term effect (Inoue et al., 2024). The Markov property is also respected as the future state (congestion levels, vehicle positions, and traffic flows) depends only on the current state and the decisions made. With these criteria being met, RL should be a fitting solution method to solve traffic congestion problems.

There are multiple cases found that apply a form of RL to optimize urban traffic control. The results are promising as waiting time is reduced (Hegde et al., 2024). Lin et al. (2018) show that a DRL algorithm can outperform traditional rule-based approaches. They suggest further research to be done regarding methods to transform traffic networks into state representations that a neural network can effectively process.

McCluskey et al. (2016) mention the potential for RL in traffic signal control. These traffic signal control cases have been researched and worked out various times applying forms of DRL that yield positive results (Li et al., 2021; Genders & Razavi, 2016; Zheng et al., 2019). The results found the delays to be reduced by 14 to 87% as well as an average queue length reduction of 17 to 73% (Rasheed et al., 2020).

A challenge when applying RL to traffic management is transforming the traffic conditions into a state space. This is difficult because the state space can become too vast to learn (Lin et al., 2018). The data collection for the agent to train on for an online learning approach in itself is also considered a challenge. A limited amount of data is available whereas a very large amount of data is required to train a DRL algorithm effectively. For offline learning via a traffic simulator, there are concerns about how accurately the reality is simulated (Han et al., 2023). Lin et al. (2018) also mention the accuracy of traffic models to be a limiting factor for the implementation of RL. However, there is active research being done on developing traffic simulation tools specifically to train DRL algorithms to optimize traffic (Wu et al., 2017; Zhang et al., 2019).

Inoue et al. (2024) mention that artificial intelligence methods using neural networks and reinforcement learning perform well for rapidly changing traffic patterns. However, the large amount of control variables and model parameters makes it computationally intense when applied to a network of more than a few traffic lights. A way to solve this is by reducing the amount of control variables. Instead of trying to optimize every local traffic light that is computationally intense, optimize a simpler control parameter such as modifying flow in a certain direction. This should also require less micro-scale information such as queues for each lane of each traffic light but rather macro-scale parameters that describe the congestion of a network. By limiting the dependent variables and simplifying the action it should be easier to learn for a reinforcement learning algorithm which in turn reduces computational time.

## 2.4   Formulation of RL approaches

Most RL approaches in the literature include a Markov decision process (MDP) of the problem. The components of the MDP are as follows: State, action, reward, and transition. The transition component is not covered here as all approaches use a traffic simulator for the state transitions. These findings are respectively shown below in Table 1, 2, and 3. Note that these formulations are for general traffic light optimization where many studies focus on a single traffic light.

| State formulations | |
|---|---|
| State | Source |
| Number of vehicles waiting at each incoming lane except the right lane and the phase of the traffic light. | X. Huang et al., 2021 |
| Current and previous speed limit as well as average speed and density of highway. | E. Walraven et al., 2016 |
| Vehicle density, speed of vehicles, direction of vehicle movement, distance between intersection and vehicle, number of edges (segments) in a lane, and last traffic signal in a lane. | Paul & Mitra, 2020 |
| Queue of each lane at a point in time at an intersection with k incoming lanes. As well as the phase of the traffic light. | J. Gu et al., 2021 |
| Queue length, phase of the traffic light, time of day, and day of week. | M. Muresan et al., 2021 |
| Incoming traffic in a given phase. As well as the number of incoming vehicles per lane. | R. Zhu et al., 2023 |
| Incoming and waiting traffic per lane. | Y. Wang et al., 2022 |

Table 1: Overview of state formulations.

It can be observed from Table 1 that the states are generally described as the number of vehicles waiting at the traffic light as well as the current phase of the traffic light. Some sources also add additional information to the state such as the speed of the vehicle or a time index. This describes the state more accurately but also increases the dimension of the state space.

| Action formulations | |
|---|---|
| Action | Source |
| Alternation and duration of traffic lights on each incoming approach. As well as the phase selection for the next period. | X. Huang et al., 2021 |
| Speed limits. *Note that this is applied for highways. | E. Walraven et al., 2016 |
| Turn a traffic light on or off. | Paul & Mitra, 2020 |
| Keep the current phase or move to the next phase. | J. Gu et al., 2021 |
| Keep the current phase or move to the next phase. | M. Muresan et al., 2021 |
| Change the times of the phase cycle. | R. Zhu et al., 2023 |
| Switching to the next phase or predefined set of phases. | Y. Wang et al., 2022 |

Table 2: Overview of action formulations.

It can be observed from Table 2 that most approaches have the changing of traffic light phases as action while respecting the constraint of a minimal green time. In contrast to most cases in the literature that use frequent actions for one traffic light, this research instead aims for infrequent actions over a network of traffic lights.

| Reward formulations | |
|---|---|
| Reward | Source |
| Negative reward for queue length per lane | X. Huang et al., 2021 |
| Number of vehicles per hour. Punishment if the number of vehicles per hour decreases. *Used for highways. | E. Walraven et al., 2016 |
| Positive/negative reward for increasing/decreasing waiting time | Paul & Mitra, 2020 |
| Number of passed vehicles divided by the number of stopped vehicles | J. Gu et al., 2021 |
| Vehicles discharged * a reward x – vehicles waiting * a reward y. | M. Muresan et al., 2021 |
| Individual reward: Number of vehicles waiting per lane per stage. Global reward: Average vehicles waiting on all lanes per phase. *Multi-agent. | R. Zhu et al., 2023 |
| Negative reward for all the combined queue lengths | Y. Wang et al., 2022 |

Table 3: Overview of reward formulations.

It can be observed from Table 3 that most sources choose queue length as the reward. Some have tried to formulate this differently by giving a reward to a ratio of vehicles that drive through the traffic light or use weights for passing and waiting vehicles. For the weights approach, the reasoning was to improve learning speed by giving an additional negative reward when vehicles wait.

## 2.5 Conclusion literature study

Traffic increases at a faster rate than the infrastructure supporting it which causes congestion. Traffic management aims to reduce congestion as well as optimize the use of road infrastructure. Techniques such as traffic light control and perimeter control help manage traffic by either guiding vehicles to alternative routes or regulating inflows to prevent congestion. Techniques such as max pressure control enable localized decisions, based on real-time data to improve traffic flow at intersections. A challenge for traffic management is the stochastic nature of traffic. First, the future traffic state must be predicted. Second, given the prediction, a solution to the traffic problem must still be found which takes long-term effects into account. A solution approach to the first challenge, predicting traffic, is supervised learning using neural networks. Neural networks allow for the generalization of data which in turn enables predicting the future state of traffic. The results are promising but the need for large amounts of data remains a challenge. A solution approach to the second challenge, finding a traffic management solution that takes long-term effect into account, is Deep Reinforcement Learning (DRL). DRL can use the strengths of neural networks and Reinforcement Learning (RL) to make sequential decisions under uncertainty with long-term effects. DRL has promising applications in traffic management as it has been shown to reduce waiting times, delays, and queue lengths. However, challenges such as scaling in large traffic networks, the high computational costs, and the accuracy of simulated training environments must be addressed for widespread implementation.

# 3 Modeling the problem and the solution method

This chapter will give additional insights into the problem and the requirements for the solution. Followed by explaining which tools are used to solve the problem. Rounded off by a solution method and the steps required to work towards that solution. This chapter will conclude by answering the following research question: "How can the solution method be designed to reduce traffic congestion while being able to provide decision support to traffic engineers".

## 3.1 Breakdown of the problem

In simple terms, the problem is that there is too much traffic within cities, which causes congestion. Traffic engineers in essence have two options: First, do their best to solve this problem reactively by diverting traffic or moderating in/out flows. A second alternative is to program a rule-based scenario that consists of many if statements. This rule-based approach will then apply preset solutions if certain traffic conditions are met. Correctly applying this to a whole city is cognitively hard, a lot of work, and difficult to oversee and maintain. A solution method for optimization problems is to mathematically compute an optimal solution for each possible scenario. However, this is not computationally feasible for traffic management due to the stochastic nature of traffic and the sequential decision-making with long-term effects. Consequently, mathematical optimization is not possible. Therefore, a different approach must be used to predict and prevent congestion in a prescriptive manner.

## 3.2 Markov decision process

As already seen in the literature, RL problem formulations often include a Markov decision process (MDP). This sub-section will go over the MDP that is used, the result can be seen below. Note that multiple variations of state space formulations, actions, and rewards have been considered in this research.

- **State:** The state contains information to represent the current state of the environment. This is a vector containing: Average lane occupancy, total vehicles, total vehicles inside the perimeter, total vehicles outside the perimeter, average speed, total queue length, average queue length, the queue outside the perimeter, the queue inside the perimeter, total travel time, average travel time, inflowing vehicles and outflowing vehicles.

The state space can be mathematically represented as follows:

**State ($S$):** The state $s \in S$ is a vector that contains the relevant traffic information in the network at time $t$:

$$s_t = \begin{bmatrix} \text{average lane occupancy,} \\ \text{total vehicles,} \\ \text{total vehicles inside perimeter,} \\ \text{total vehicles outside perimeter,} \\ \text{average speed,} \\ \text{total queue length,} \\ \text{average queue length,} \\ \text{queue outside perimeter,} \\ \text{queue inside perimeter,} \\ \text{total travel time,} \\ \text{average travel time,} \\ \text{inflowing vehicles,} \\ \text{out flowing vehicles} \end{bmatrix}$$

Each element in $s_t$ is a scalar that represents a specific traffic feature at time $t$. This should provide the required information for the neural network to understand the environment.

Throughout this research addition state space formulations are attempted to find the best formulation. The additional state space formulations that are considered include adding more information such as queues for each lane as well as less information by using only a part of the aforementioned state space. The state space seen above is most the successful. Additional details on state space experimentation can be found in Section 4.3.

- **Action:** The action is a vector containing scalars with which the respective phase duration will be multiplied. For example, if the first element is equal to 0.75 the original green time of the first phase (28 seconds) will now become 28*0.75 = 21 seconds.

This can be mathematically formulated as follows:

**Action ($A$):** The action $a \in A$ is a vector that adjusts the green phase durations for traffic lights in the network. Let $a = [a_1, a_2, a_3, a_4]$, for each $a_i$, with $i = [1, 2, 3, 4]$, be a multiplier for the duration of the original phase $d$, with $d = [28, 11, 28, 11]$. The new phase duration $c$ is then computed as follows:

$$c = d \times a = [d_1 \times a_1, d_2 \times a_2, d_3 \times a_3, d_4 \times a_4]$$

Similar to the state formulation, multiple action formulations are attempted throughout this research. Simpler versions such as only adjusting the main phases' duration or choosing a ratio between the main phases have been tried. The action formulation above is the most intuitive formulation and yields similar results to the other formulations.

- **Reward:** The reward is the mean speed over the whole network.

The reward can be mathematically formulated as follows:

**Reward (R):** The reward $R(s, a)$ for a state $s$ and action $a$ is the mean speed of all vehicles in the network, which can be represented as:

$$R(s_t, a_t) = \frac{1}{N} \sum_{i=1}^{N} v_i$$

Where $N$ is the total number of vehicles in the network at time $t$, and $v_i$ is the speed of vehicle $i$.

This research also experimented with various reward formulations, using the mean waiting time, trip completion rate, and weighted combinations of the average speed, mean waiting time, and trip completion rate. The reward above is the best-performing reward and also the most intuitive.

- **State transitions:** The state transitions will be handled by SUMO (Simulation of Urban Mobility), a traffic simulator. SUMO will perform the steps in the simulation which moves all the cars to the next location as well as controlling the traffic lights. After many steps, the next state is forwarded to Python. Here the DRL algorithm chooses an action based on the state. This action is then sent to SUMO, which in turn will adjust the phase durations of the traffic lights on the perimeter accordingly.

The state transitions can be mathematically formulated as follows:

**State Transitions ($P(s'|s, a)$):** The state transitions are managed by SUMO, which updates the environment. Formally, the transition probability $P(s'|s, a)$ represents the probability of moving to a new state $s'$ given a current state $s$ and action $a$. Since the exact transition probability is complex and done by SUMO, it is described as:

$$s_{t+1} = f(s_t, a_t)$$

Where $f$ is a function (simulated by SUMO) that takes the current state $s_t$ and action $a_t$ and outputs the next state $s_{t+1}$ based on vehicle movements, traffic light changes, and other dynamics in the simulation.

- **Stages:** The duration between actions is kept at 90 seconds, the time of a traffic light cycle.

The goal of this research is to provide decision support at a larger time scale of 15 minutes, however, this is not feasible due to computational reasons. Only making a decision every 15 minutes would reduce the samples per iteration by $(15 * 60)/90 = 10$. As a result, training the reinforcement learning algorithm will take roughly 10 times as long. So for simplicity, 90 seconds is taken, if this approach works it is worth experimenting with a stage of 900 seconds.

## 3.3 Requirements of the solution method

- Technolution puts significant emphasis on wanting to give decision support, thereby augmenting the intelligence of the tool's user. This means that the output solution must be easy to interpret for human users.

- The solution should not have to be applied every couple of seconds, but rather over a larger time interval. The results from the literature research in Section 2.4 show that most approaches focus on optimizing a single traffic light. The results would then output a phase selection every 5 seconds. Although this approach is proven to be effective the results might not be easy to grasp. In addition, an output is given in a very short time interval. This makes it not suitable for decision support as a traffic engineer will not input direct traffic phase durations. Therefore, a macro-scale approach is chosen to optimize a network, instead of looking at the micro-scale optimization of a single traffic light. The goal here is to output an intuitive solution over a longer time horizon.

- The solution method must also be able to deal with the complex nature of traffic management. As well as finding an alternative to the mathematical approach of iterating the solution for every possible traffic situation. Since this is not feasible with the size and stochastic nature of the problem.

In summary, the solution method should be easy to interpret, over a long time interval and find an alternative to mathematical computation for every situation.

## 3.4 Solution method to solve the problem

This subsection will first discuss the tools that are used, followed by the solution method to solve the two challenges of traffic management found in the literature. The first challenge is to predict traffic due to its stochastic nature. The second is finding actions, which solve the traffic problem that take into account long-term effects. Neural networks are chosen to solve the first challenge, predicting traffic, for the following abilities: To learn complex problems, being able to handle large state spaces resembling traffic, adaptability to real-time data, and the ability to make fast predictions once trained. Deep Reinforcement Learning (DRL) is chosen as the solution method to the second challenge, finding actions that solve traffic problems that take into account long-term effects. DRL is chosen for the following abilities: Scalable to large networks using neural networks to approximate a policy as well as the ability to make decisions under uncertainty while taking long-term effects into account. In addition to the ability to continuously train and improve on new data and the ease of using a simulation tool to train in. The steps in which these tools are used are described in the next paragraph.

The idea is to use DRL to find the best action to prevent future congestion. The output must be intuitive to understand and be applied over a larger time period, given the requirements stated in Section 3.3. As previously mentioned, this eliminates single-traffic light optimization. A good solution method for a larger area is perimeter control. Limiting the inflow to an area so that the traffic remains below a certain threshold, can be seen in Figure 3. The red line marks the protected region on which perimeter control will take place. Along the red line, the flow into and out of the city will be moderated. Given that in this case a neural network can predict the future traffic inside

a city with sufficient reliability, a DRL algorithm can find the optimal action for a traffic state to achieve the highest reward. The reward that is chosen to be optimized is the mean speed across the entire traffic network, which will be maximized. The DRL algorithm should then give as output an inflow/outflow moderation. This will be a vector that has an element for each phase (excluding yellow). The elements in this vector represent a percentage increase or decrease per phase. This is intuitive to understand as it simply tells a traffic manager per direction if the durations should be increased/decreased. In addition, the action of modifying phase durations would not need to be applied every 5 seconds like local traffic light optimization. This is because the inflow/outflow of vehicles over 5 seconds will not significantly impact the average flow of the whole city network. The approach of perimeter control thus covers all requirements as the output is intuitive and over a longer time period.

The problem that does arise is that local optimization of traffic lights is still required. Simply moderating the duration of the inflowing phase for every single traffic light on the perimeter by a fixed amount is not optimal. In reality, traffic is not uniformly distributed at every entry point of the network. The result of simply increasing green time uniformly for all traffic lights in a direction would be wasted green time. To solve this problem and prevent increasing green time unnecessarily, a second control layer is added: Max pressure (MP). MP will regulate the traffic lights at a local level to prevent a waste of green time at local traffic lights. Meanwhile, perimeter control will regulate the congestion within the city on a global level by adjusting in and out flows at the perimeter. While perimeter control prevents congestion by increasing or decreasing global flows, MP makes sure that these flow modifications are being used in the most efficient way possible at a local level per traffic light. These methods are then combined to have both local and global optimization of traffic. See Section 2.1 for more detailed information on the working of max pressure and perimeter control. This approach will give an understandable output to the traffic manager every 15 minutes.

To summarize the main challenges and requirements with the respective design choices for the solution method:

- Since traffic is complex, stochastic, and changes have long-term effects, the solution method must be able to deal with this challenge. DRL is chosen as the tool to solve the problem as it combines the strengths of neural networks and reinforcement learning to deal with sequential decision-making under uncertainty in complex state spaces.

- Since the output is used for decision support, it must be easy to understand. To achieve this, the percentage adjustment per phase is provided as an action, which should be intuitive and straightforward for traffic engineers.

- The output should also be over a longer period of time to prevent the traffic engineer from having to make constant changes. To achieve this, a macro optimization approach is chosen due to the longer time periods between actions. The specific method is perimeter control, adjusting the flow from and to the city for the traffic lights on the perimeter.

- Since perimeter control uniformly adjusts traffic at a global level and traffic is not uniform, a form of local optimization is required to prevent wasted green time. This will be done by max pressure to ensure that the adjusted in and outflows are well distributed.

- The goal in traffic is to have everything moving as fast as safely possible. A good measure for this is the mean average speed in the whole network. The mean speed is also measurable in the real world. For these reasons, this has been chosen as the optimization parameter.



Figure 3: Illustration perimeter control.
Perimeter control and thus adjustment of traffic light phase durations will take place at traffic lights on the red line. Note that the grid network used in this research is 5x5, not 7x7 as seen in the figure (Kampitakis & Vlahogianni, 2024).

## 3.5 Breakdown of the solution steps

Before implementing the final solution method of DRL with perimeter control and max pressure intermediate steps should be taken. The steps taken to work towards the problem are as follows:

First, supervised learning can be used to predict traffic congestion. In this case, a neural network can be trained to predict a label (future traffic condition) given an input (current traffic condition). In contrast to a mathematical approach, a neural network does not require data for every single traffic possibility. Instead, a neural network is designed to generalize the data by recognizing

patterns. This is done for various states to determine the optimal state formulation. Successfully training the neural network to make predictions covers the first part of the problem: Predicting upcoming congestion.

Second, after it is possible to predict congestion, a preventive measure for future congestion should be found. Optimizing sequential actions based on predicted outcomes from a neural network falls under the domain of Deep Reinforcement Learning (DRL). The goal of RL is to find the best action given a certain state by maximizing the reward. DRL in turn is a form of RL that uses a neural network to predict the best action, given a state. This replaces the need for a traditional tabular approach or function approximators. Now, also the second part of the problem is covered: Finding solutions to prevent upcoming congestion.

Third, to verify that everything works as intended a simple action for the DRL algorithm is chosen to begin with. This approach will naively modify the duration of the traffic light phases uniformly for all traffic lights. This action is the action that is presented in Section 3.2. This solution method is not expected to perform very well due to the lack of local optimization but should outperform the static baseline policy.

Fourth, the DRL algorithm will use perimeter control in combination with local optimization. Note that this local optimization is not yet done by the final local optimizer, max pressure, but instead, a different integer linear programming formulation that has been proven to work in combination with DRL (Kampitakis & Vlahogianni, 2024). This is done to again verify that everything is working as intended while working towards the final goal. In this approach, the action made by the DRL algorithm is the amount of cars that are allowed to enter the perimeter during the next time step. The amount of cars is then divided among the intersections on the perimeter depending on the queue present at the intersection. This local optimization is done by solving a set of linear equations.

Lastly, perimeter control is combined with max pressure to optimize traffic at a global and local scale. An in-depth explanation of these two techniques is given in the following two sections as well as explaining how and why both techniques should be combined.

## 3.6   Max pressure

This section covers how max pressure (MP) can optimize local traffic lights using an optimizer. A complete mathematical notation and an integer linear problem formulation are presented. The max pressure approach uses pressure to describe how much green time a phase requires. Pressure is calculated as follows: The difference between the queue length of the incoming lane and the outgoing lane. A more in-depth explanation can be found in Section 2.1.

### 3.6.1   Notations

An overview of the notations, adapted from Tsitsokas et al. (2022), can be found below:
N = The set of intersections, elements n
Z = The set of links, elements z
I = The set of incoming links, elements i
O = The set of outgoing links, elements w

D = The set of directions, elements d
J = Phase of the traffic light
$C_n$ = Cycle time
$C_z/C_w$ = Capacity
S = Saturation flow
B = Turning movement rates
G = Green time
k = Control cycle
L = Fixed total loss time

### 3.6.2 Max pressure formulation

This section covers how MP can be described mathematically and implemented. A mathematical description of the max pressure policy along with the constraints is shown below. These formulas are adapted from Tsitsokas et al. (2022). The constraints are as follows:

$$\sum_{j \in \mathcal{F}_n} g_{n,j}(k_n) + L_n = C_n \tag{1}$$

$$g_{n,j}(k_n) \leq g_{n,j,min} \quad j \in F_n \tag{2}$$

Formula 1 prevents the sum of green durations + lost time (yellow) from exceeding the cycle time. Formula 2 makes sure that the green duration chosen is larger than the minimum green time required. This minimum green time is mandatory for all traffic lights, see (Handboek Verkeerslicht-enregelingen 2022, 2022). The amount of pressure $p_z(k)$ of every incoming link $z \in I_n$ of node n for control cycle k is calculated as follows:

$$p_{z,d}(k) = \left[ \frac{x_{z,d}(k)}{x_{z,d,max}} - \sum_{d \in \mathcal{D}} \frac{x_{w,d}(k)}{x_{d,max}} \right] S_z, \quad z \in I_n \tag{3}$$

Formula 3 calculates the pressure as the queue length of the current link divided by the link's capacity - the sum of the pressure for the lanes connected to the link. It is important to note that here a lane is a part of the road, not the complete road. Each lane has its own traffic light and pressure. Since negative pressure is meaningless an additional constraint, Constraint 4 is added:

$$p_{z,d}(k) \geq 0 \tag{4}$$

### 3.6.3 Phase selection

For phase selection two approaches can be used: freedom of phase selection or a predefined order of phases. The first option can be applied as follows: A phase gives green to a certain direction at the traffic light. Some phases can be active at the same time. To choose which phases to activate, the total pressure of the combined possible phases is determined. For example: phases 1 and 5 can have green simultaneously. The total pressure would be the pressure of the south lane in direction left + the pressure of the north lane in direction left. An overview of the visual representation of the

phases can be found in Figure 1. The pressure is calculated for all possible combinations of phases. The second option is simpler. It is only necessary to calculate the total time per phase, based on the modeling assumption that the phases follow a fixed order. This is combined with the assumption that every phase will be activated every cycle with a minimal green time. This simplified approach will be considered in this study.

### 3.6.4 Determining the total cycle length

Determining the green time duration can be done in two ways. The first approach is to set a fixed cycle length where the time is filled up by having each phase active for a certain amount of time. A second possibility is to apply a cycle-free approach where there is no set time in which a cycle of phases must be completed. This creates a lot of freedom but also difficulties. The problem with a cycle-free approach is how to determine the duration of the phase. Keeping a phase for longer allows more traffic to pass but will have other lanes wait more as a direct consequence. Frequent switching of phases, due to a short cycle length, leads to more lost time due to the yellow time. This is a difficult trade-off to optimize. Instead, the simpler approach of a fixed cycle is chosen. Now the challenge is to determine the total cycle time. The method to determine the total cycle time is the Webster method adapted from Yadav et al. (2022). The following notations are used:
N = The set of phases, elements n
$C_0$ = Total cycle time
L = Lost time per cycle
R = All red time
y = Flow ratio
q = Normal flow
s = Saturation flow rate

This method determines the total cycle time as follows:

$$C_0 = \frac{(1.5L + 5)}{(1 - Y)} \tag{5}$$

With

$$L = 2n + R \tag{6}$$

$$Y = y_1 + y_2 + ... + y_n \tag{7}$$

$$y_n = \frac{q_n}{s_n} \quad \forall n \in N \tag{8}$$

The saturation flow rates are taken from Yadav et al. (2022). These have the following values per width of the lane:

The value for q can be obtained via data analysis. Substituting the values for q and s in Equation 8 and adding these for all phases gives the flow ratio. After filling in Equation 6 to get the lost cycle time, all variables for Equation 5 are known. Filling this Equation in returns the total cycle time.

| Width in meters | 3.0 | 3.5 | 4.0 | 4.5 | 5.0 | 5.5 |
|---|---|---|---|---|---|---|
| Saturation Flow (PCU/hr) | 1850 | 1890 | 1950 | 2250 | 2550 | 2990 |

Table 4: Standard values for saturation flow, S, for widths 3 to 5.5 meters according to Webster.

### 3.6.5 Distribution of green time

The distribution of green time can be determined by solving an optimization problem based on a set of equations describing the goal and the constraints. These steps and equations are adapted from Varaiya, P. (2013).

First the green time according to MP is determined. The total green time is the cycle time minus the lost time. This can be described as follows:

$$G_n = C_n - L_n \tag{9}$$

Second, the green time per phase is defined as follows:

$$\tilde{g}_{n,j}(k) = \frac{P_j(k)}{\sum_{i \in F_n} P_i(k)} G_n, \quad j \in F_n \tag{10}$$

Here the pressure per phase is divided by the total pressure to get the proportion of pressure as a result. This is multiplied by the total green time to get the total green time per phase. To ensure that the green time per phase found in Equation 10 also complies with the traffic light guidelines, an additional optimization problem is formulated:

$$
\begin{aligned}
\underset{G_{n,j}(k)}{\text{minimize}} \quad & \sum_{j \in F_n} \left( \tilde{g}_{n,j} - G_{n,j} \right)^2 \\
\text{subject to} \quad & \sum_{j \in F_n} G_{n,j} + L_n = C_n, \\
& G_{n,j} \geq g_{n,j,\min}, \quad \forall j \in F_n, \\
& |G_{n,j} - g_{n,j}^p| \leq g_{n,j}^R, \\
& G_{n,j} \in \mathbb{Z}^+.
\end{aligned}
\tag{11}
$$

The aim is to have the allowed green time $G_{n,j}$ be as close as possible to the green time per phase found by MP $\tilde{g}_{n,j}$. The constraints are as follows: First, total green time + lost time is equal to total cycle time. This is the same as Formula 1. Second, the green time must be larger than the minimal green time. Third, there is a maximum change of green time $g_{n,j}^R$ allowed between the current cycle $G_{n,j}$ and the previous cycle $g_{n,j}^p$. This optimization problem can be solved by for example Optuna optimizer.

Despite the requirement to give decision support, it is not possible to have a large time step for micro-level traffic control, in contrast to macro-level control. On a micro-level, the time step chosen will be 5 seconds. This means the selected phase will be revised every 5 seconds. During this

revision, the pressure for all lanes will be calculated again and a phase will be selected accordingly. The part that will be used for decision support is perimeter control which will be used on a macro scale with a significantly larger time step of 15 minutes. This is covered in the next subsection.

## 3.7 Perimeter control

The goal of perimeter control is to reduce congestion in the perimeter by changing the green times of the traffic lights on the perimeter. This is done by adjusting the weights for the pressure that the MP method uses to give green durations, the weights are given in Equation 12. Here the weights modify the amount of pressure a lane has, which then indirectly modifies the in and out flow of that lane. To facilitate the addition of weights, the formula for pressure 3 is modified to include these weights $w_d$:

$$p_{z,d}(k) = w_d \left[ \frac{x_{z,d}(k)}{x_{z,d,max}} - \sum_{d \in \mathcal{D}} \frac{x_{w,d}(k)}{x_{d,max}} \right] S_z, \quad z \in I_n \tag{12}$$

When analyzing Equation 12 it can be observed that the pressure is just the difference between upstream and downstream occupancy of the links. By adding a weight to this pressure, it is possible to adjust the interpretation of the pressure by the controller. If a factor smaller than one is used the perceived pressure is less and a shorter green time duration will be provided. This way the traffic into the perimeter is moderated. The opposite holds for weights greater than one. Intersections not on the perimeter will have a weight of one which turns Equation 12 back into Equation 3.

## 3.8 Supervised learning using neural networks

Supervised learning is a subgroup of machine learning where models are trained on labeled data to learn the relationship between the input features and the target outputs. In this research, it will be applied to predict the future traffic situation based on the current traffic situation. The labeled data is the current traffic situation and the target output is the future traffic situation, in this case, the mean speed. Different state formulations will be tried to find the best state formulation to describe the traffic situation. By using neural networks it is possible to predict the future traffic situation based on the input of the current traffic situation.

## 3.9 Iterative supervised learning

As described previously, supervised learning can predict a target or label, given a set of input features, which can be used to make decisions. If the traffic light configuration is part of the state vector that describes the current traffic situation, it is possible to predict the best traffic light configurations based on the current traffic situation. Applying this decision-making sequentially to a traffic simulation should improve the mean speed. Doing this process iteratively will continuously improve the mean speed as well as the accuracy of the predictions when using the newfound data to fine-tune the existing neural network. After multiple iterations, the mean speed improvements should decrease and eventually converge. The flowchart of this process can be seen in Figure 4. In conclusion, repeatedly using a neural network to make decisions based on the current state of traffic can be a way of improving the mean speed and thus reducing congestion.

Figure 4: Flowchart for iterative supervised learning approach.

## 3.10    Combination of perimeter control and max pressure using DRL

After making decisions using a neural network, the next step is to use Deep Reinforcement Learning (DRL) to make sequential decisions under the uncertainty of traffic. The DRL algorithm can be used to output weights that modify the pressure of a lane and in turn the in and out flows. As described earlier, a second layer of control is required for local optimization. This will be done by solving Equation 12 for all lanes at each traffic light. The result of using DRL to formulate weights for the pressures followed by max pressure optimization should lead to global and local optimization of traffic lights.

## 3.11    Choice of DRL algorithm

The action space of the problem should be a continuous range of values to modify the in/out flow and not be limited to a fixed set of options. As the action space of the problem is continuous, discrete action space RL algorithms are eliminated as a possibility. One of the most popular continuous action space DRL algorithms is proximal policy optimization (PPO) (Schulman et al., 2017). The reason for this popularity is due to the ease of implementation whilst having relatively high performance (Van Heeswijk, 2023). Due to the ease of implementation and allowing a continuous

action space PPO is chosen as the DRL algorithm. In addition, twin-delayed deep deterministic policy gradient(TD3) is applied as well. The PPO and TD3 implementations are taken from the StableBaselines3 library (Raffin et al., 2019).

## 3.12    Concluding approach for modeling the problem and the solution method

The difficulties that are found in solving traffic congestion problems is the ability to effectively predict and solve them. A solution for this is using a neural network to predict the future traffic state and in turn congestion. By using deep reinforcement learning (DRL) the best action can then be found given a traffic state. The DRL problem formulation can be described as follows. The current traffic conditions in the network are used as an input state. The in- and outflow moderation as action. The mean speed throughout the network as the reward. The traffic simulator SUMO will be used for the state transitions. The stage duration will be 90 seconds.

To make sure that the solution provided by the DRL algorithm is over a larger time period, global optimization of the city network is done, as global optimization uses larger time intervals between actions. The global optimization method used is perimeter control. To prevent wasting green time on the boundary of the perimeter, perimeter control is combined with max pressure. See Figure 5 for a flow chart showcasing how the phase durations are determined and Figure 6 for the flowchart showing how the process is done during simulation. Max pressure will optimize local traffic lights while perimeter control optimizes global traffic lights. The DRL algorithm will regulate in- and outflow of the perimeter. This solution is both intuitive and operates over a longer time scale, aligning well with the solution requirements.



Figure 5: Flow chart displaying process for determining phase durations.
The DRL algorithm does PC by modifying the weight of queues which will change how heavily each lane contributes to the phase duration in the calculation for MP, see Equation 12.

Figure 6: Flow chart displaying the decision-making during the simulation.
Showcasing the constant use of MP at the local level and the occasional modification of queue
weights by PC determined by the DRL algorithm.

# 4 Experimental settings

This chapter covers the experimental settings and what choices are made. The development of the solution method is done step-wise and starts with training a neural network on traffic data to predict future traffic. This neural network is then used to choose the best phase durations during a simulation to represent DRL in a simple way. Afterward, DRL is used with a form of naive perimeter control. Here the phase durations of the traffic lights on the perimeter are all changed equally without the use of local optimization. The before last step is to recreate the perimeter control approach from Kampitakis & Vlahogianni (2024) with a local optimizer, different from max pressure. The final step is to combine perimeter control with max pressure. These intermediate steps are done to slowly increase complexity while maintaining insights into the changes and possible problems. The chapter concludes by answering the following research questions: "How should the RL algorithm be applied to optimize traffic traffic flow?".

## 4.1 Modeling choices and assumptions

For simplicity, the traffic network is modeled as a square grid network. A grid network might not be as common in Europe but these are commonly found in America and thus represent a real-world traffic network. The network consists of 25 traffic lights connected by 200-meter double-lane two-way roads. These network settings are chosen to match the network used by Kampitakis and Vlahogianni (2024), except a bit smaller using a 5x5 network instead of a 7x7. The network is created in Simulation of Urban Mobility (SUMO), a commonly used traffic simulator. The network can be seen in Figure 7. Python will be used to interact with SUMO.

Traffic will be created to arrive and depart throughout the entire network to represent a real city. Traffic will follow a thousand randomly generated routes. However, this process is not entirely random. To prevent certain roads from being disproportionately included in a significantly higher number of routes, and thereby experiencing increased traffic, the following solution is implemented: Giving probability weights that decrease the probability of selecting frequently used roads in the random route-generating process. The result is that the likelihood of roads being disproportionally included is significantly reduced. This makes the traffic a bit more uniform throughout the city without being exactly evenly distributed. Roads are on average included in 30-40 routes. This will make sure that all parts of the network are used. The majority of the traffic will be from the outside of the perimeter to the inside of the perimeter. This is chosen to represent a morning rush where people from outside the city go to work in the city. The arrival process along these created routes is a Poisson arrival process. When the simulation is started a warm-up state is loaded, this is the same every single time. This is computationally faster than running the traffic simulation for a couple of minutes for the network to fill up and also makes the starting states more realistic. To prevent overfitting, the routes are randomly generated for each iteration. Combined with the stochastic arrival rate this should make the traffic sufficiently random.

The intensity of the created traffic is between the following three conditions: Low traffic, almost no waiting times, and everything flowing well. Medium traffic, congestion starts to occur and queue times become longer. High traffic, the network is becoming fully congested and queues become very long. Note that traffic not only consists of these 3 types but also many in-between types of traffic.

For configuring the traffic light settings the following assumptions are made. To make the problem realistic there is a minimal green time. This is implemented as a constraint, preventing the algorithm from turning the light off within 8 seconds. There also is a minimal time between two conflicting directions switching from red to green. This is covered by the in-between yellow phase of 3 seconds. These are common parameters in the real world. The optimization of local traffic lights using max pressure will be done according to the formulas provided in Section 3.6. Pyomo, an open-source optimization package, will solve the set of linear equations.

The states are chosen in a way that these can be applied in a real-world scenario and thus must also be measurable. Because of this for example traffic density is excluded as this is only something that can be accurately measured in a simulation. According to Edwin Mein (personal communication, July 25, 2024), an expert in the field of traffic management, a more realistic approach is to use average traffic speed. This data is already available and would not require any new infrastructure to be built. This also gives a good indication of the congestion in a city as speed and congestion are directly related. Alongside average speed, additional state features are chosen, and an in-depth explanation of the state formulation is provided in the next section.



Figure 7: SUMO traffic network.

## 4.2   Data collection for neural network training

To generate data for training the neural network many simulations are done. Each simulation changes either the traffic intensity or the traffic light settings. By choosing a wide range of traffic intensities and traffic light settings a variety of data is collected. A data point is generated each traffic light cycle (90 seconds). The collected data is then saved in a text file containing the state of the traffic, the traffic intensity, traffic light settings, and the cumulative mean speed over the next five cycles. As mentioned before, for each simulation run a new set of random routes is created to prevent overfitting.

The training data does not cover the traffic intensities and the traffic light settings of the validation data. Some traffic light settings and traffic intensities for the validation data are also out of the range of the training data. With these differences in the data, it is assumed that the two data sets are significantly different and the neural network would indeed need to predict unknown scenarios. The reason for this two-step data generation approach where different settings are used for validation and training is that otherwise predicting a random part of the training data in the validation process is too easy as the data is sequential and very similar per step. If a random point separated from the training data would be selected for validation the neural network would know the state and reward just before that and just after that data point making the prediction a simple linear regression instead of testing whether the neural network understands the data and can predict new situations.

To train the neural network 12000 samples are taken. Ideally, more samples would be used but creating 12000 samples takes more than 24 hours. The validation data contains 1,500 samples. 10% of the training data is randomly selected and separated to use as validation during the training process to monitor performance. This is not to be confused with the final validation data which is kept separate. But rather the data used to update the weights of the neural network during training. The results found using the best hyperparameter and neural network architecture for different state spaces can be found in Section 5.2.

## 4.3   Training a neural network to predict future traffic situations

The first step that is taken before trying to apply a form of DRL is to check whether a neural network can predict the relations between input data (traffic state, traffic intensity, and phase durations) and labels (mean speed, reward) with sufficient accuracy. The traffic state contains information about the traffic in the network, the traffic intensity is the expected value for the Poisson arrival rate, and the phase duration is a vector of length four, showing the green time per phase. The label (output) is the sum of the mean speed over the next five cycles. To train a neural network to predict future traffic the following steps are done: First, many different traffic intensities are simulated using a wide variety of traffic light settings. Second, the data is then normalized to ensure all state features are between values of 1 and -1. The reward is also normalized to be between 1 and -1 using a different normalization vector. The normalization vector is saved to later be able to transform the data back to their real values for visualization. Third, the normalized data is then used to train a neural network. Lastly, the trained neural network is compared to the validation data to analyze the performance of the trained neural network. For the neural network training the scikit-learn library (Pedregosa et al., 2011) is used.

To find the best hyperparameter settings and neural network architecture the GridSearchCv and RandomizedSearchCV functions from scikit-learn are used. These functions try a variety of hyperparameter and neural network architecture combinations. Judging based on the mean square error (MSE) between the target and expected label, the best-performing combination is used. The hyperparameters can be seen in Table 5.

Table 5: Hyper parameters for Neural Network Training

| Hyper parameter | Value |
|---|---|
| Learning Rate | 0.001 |
| Batch Size | 64 |
| Optimizer | Adam |
| Hidden layers | 128, 64, 32, 16 |
| Activation Function | ReLU |
| Learning rate | Constant |

To determine what formulation of state space is most effective, a variety of state space formulations are used. The different state formulations that are used are shown below. Note that the current phase durations and traffic intensity are included in all states and for simplicity will not be shown in the overview of different state formulations in Table 7.

| Feature | State 1 | State 2 | State 3 | State 4 | State 5 | State 6 |
|---|---|---|---|---|---|---|
| Total lane occupancy | ✓ | ✓ | ✓ | | | |
| Average lane occupancy | ✓ | | ✓ | ✓ | | |
| Total vehicles | ✓ | ✓ | ✓ | ✓ | | |
| Total vehicles inside perimeter | ✓ | ✓ | ✓ | ✓ | | ✓ |
| Total vehicles outside perimeter | ✓ | ✓ | ✓ | ✓ | | ✓ |
| Average speed | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Total queue length | ✓ | ✓ | ✓ | | | |
| Average queue length | ✓ | | ✓ | ✓ | | |
| Queue outside perimeter | ✓ | ✓ | ✓ | ✓ | ✓ | |
| Queue inside perimeter | ✓ | ✓ | ✓ | ✓ | ✓ | |
| Total travel time | ✓ | ✓ | ✓ | | | |
| Average travel time | ✓ | | ✓ | ✓ | | |
| Inflowing vehicles | ✓ | ✓ | ✓ | ✓ | | |
| Outflowing vehicles | ✓ | ✓ | ✓ | ✓ | | |
| Queue length for every traffic light | | | ✓ | | | |

Table 6: Comparison of state formulations and their features.

## 4.4    Iterative supervised learning

Before applying reinforcement learning, the trained neural network will be iteratively fine-tuned. A simulation of traffic is done as usual but now every five cycles (450 seconds) the neural network predicts the best traffic light phase durations for the upcoming five cycles. This is done as follows: First, the current traffic situation data from the simulation is gathered and put into a state of the same form that the neural network has been trained on. Second, for a large list of phase durations the label (mean speed over the upcoming five cycles) is predicted. Third, the phase durations that give the highest mean speed are selected and applied in the traffic simulation for the next five cycles. The process is then repeated until the simulation ends. This is done for a large amount of traffic intensities and various starting traffic light configurations. The data created from these simulations is then used to fine-tune the existing neural network. While the training data is sufficient to show the increase in mean speed over iterations, an additional validation run is required to show the new found prediction accuracy of the fine tuned neural network. To do this the same process as for training is repeated to create validation data but then with different, mostly unseen, traffic intensities. A flowchart of the iterative supervised learning process can be seen in Figure 4. Performing iterations of supervised learning should increase accuracy in predicting future traffic. Additionally, the mean speed of all simulations should be higher for the data where the neural network makes decisions. The mean speed should improve per iteration but the degree of improvement should decay and eventually converge. This approach should yield a similar result to reinforcement learning but is easier to understand. For this reason, the iterative supervised learning approach is done before DRL. The results can be found in Section 5.3.

## 4.5    Naive perimeter control

To work towards the final solution, individual components are first created and tested, as it is incredibly difficult to debug the final model of combining max pressure and perimeter control. First, perimeter control is used on the outer traffic lights excluding the corner points. All traffic lights in the network have standard phase durations and a fixed cycle time. The RL algorithm adjusts the phase durations of the traffic lights on the perimeter. The phases of the traffic light are shown in Figure 8. These adjustments still need to adhere to the fixed cycle time. This is done by dividing the sum of the chosen durations by the cycle time. Afterwards, the action is multiplied by this number. The action is once again continuous, therefore PPO and TD3 are both applied here. The Markov decision process for this approach is discussed in Section 3.2.

A slightly modified version of the above-described method for naive perimeter control is also attempted. This modified version is created to get a better understanding of how a different action method would affect the performance. The previously mentioned version of naive perimeter control uses a vector action to modify the duration of all phases. The alternative version instead uses a scaler as action, which determines the ratio of the main phases that are parallel and perpendicular to the perimeter while keeping the secondary phases constant. Figure 8 shows all phases, here the phase that is green, starting at roads 0,1 & 6,7, is the main phase. The left turn on roads 2 and 8 are secondary phases. This approach would thus determine the ratio of green time between the main phases, roads 0,1 & 6,7 and roads 3,4 & 9,10. This is done by taking the cycle time, removing yellow time, and secondary phase durations. This gives the remaining phase duration for the two main phases. The scalar is then used to divide this time between the two main phases. The state and

Figure 8: Traffic light phases.

reward are kept the same as in the regular naive perimeter control.

## 4.6 Combining perimeter control with an optimizer

During this research, another research was being conducted with a similar idea. Kampitakis & Vlahogianni (2024) also attempt to implement perimeter control to reduce congestion inside the city. The difference in their solution is the solving of wasted green time along the perimeter as well as a different action. The action that the DRL algorithm does is determine the number of cars that are allowed to enter the network. The local optimization is then done by distributing these cars among the traffic lights on the perimeter by assigning green durations. The result is that the overall idea of combining local and global optimization is similar but the execution is slightly different. As this approach is similar, a slightly simplified version of the approach is applied as an intermediate step towards the final solution which uses DRL to combine perimeter control and max pressure. This will help to check if it is correctly implemented as well as give additional insights.

## 4.7 Final version, perimeter control with max pressure

This is the final step that is worked towards. Here a DRL algorithm uses perimeter control to moderate in/out flows which are applied to local traffic lights on the perimeter using max pressure. An in-depth explanation of this approach is given in Section 3.6 & 3.7.

## 4.8   Model verification

Verification is checking whether a model performs as it is intended to perform. In the case of this research, it is if the simulation is running as intended, are the traffic lights using the correct durations, if the traffic is generated and flowing as intended and are the state features calculated correctly.

The environment is verified using the graphical user interface (GUI) option from SUMO to observe how the traffic flows, as well as the working of the traffic lights. The code is tested using various print statements to verify that the state, action, and reward are correctly calculated and applied. The DRL algorithm is also run with an action space consisting of one action, the static standard traffic light setting, to see if it gives the same results as the simulation without the DRL algorithm. All of these tests were successful. The progress of the DRL algorithm is analyzed using Tensorboard. Furthermore, it is assumed that the DRL algorithm functions correctly and robustly unless proven otherwise since they are taken from the widely used Python library StableBaselines3.

For supervised learning, the scikit-learn library is used to train the neural network. This is a commonly used package in machine learning and is assumed to work as intended. When applying the iterative approach the action selection method is thoroughly tested to make sure that it performs as intended. This also passed all the verification tests.

## 4.9   Computational time

Training an RL agent requires a lot of data, thus Computational time is a huge bottleneck for this current research. For the DRL approach to achieve decent performance, the simulation must run for 12+ hours. The reason for this is that there are many cars, as the network is large and the traffic should represent a busy city. Simulations take longer as the cumulative amount of cars in the simulation increases. The program needs to keep track of the speed and location of all these cars causing the simulation to start taking very long as the network fills up. A simulation of an hour of traffic will take multiple minutes. This is a huge challenge faced during this research.

There are two main factors for the computational time: First, the most significant factor, simulating a congested city. As there are so many cars in the simulation the simulation speed becomes very slow. Second, for the DRL algorithm to function a lot of information must be extracted from the simulation. As the number of cars increases, so does the data that needs to be extracted, which makes it slow. The decision-making of the DRL algorithm is not the bottleneck but the simulation itself.

To account for the aforementioned challenge of computation time, computationally faster coding methods were implemented as well as running four environments in parallel to train one neural network. This allowed for faster data generation and collection. A remote desktop was also used to double the computational resources available. However, this was still not within the desired time and thereby remains a large issue.

## 4.10   Conclusion of the solution approach

In this chapter, multiple small steps are presented in working towards a final solution. These steps include first using supervised learning to predict future traffic. Second, using iterative supervised

learning to gain insights into good policies under various traffic situations. Third, using the 5x5 grid network with a simple action using DRL. Fourth, creating an approach that is proven to work using perimeter control. Lastly, implementing the final model combining perimeter control and max pressure. Applying these steps ensures gradual build-up and limits mistakes. The DRL algorithms of choice are PPO and TD3 for their continuous action space and ease of implementation. For verification, the SUMO graphical user interface will be used to see if the traffic lights work as intended. To limit computational time code optimization is done as well as parallel processing. By using parallel processing, 4 simulations were able to be run simultaneously to train one neural network.

# 5 Experimental results

This chapter covers the validation as well as the results of the approaches described in Chapter 4. This includes the following: First, all approaches are validated. Second, the performance of the neural network in terms of predicting the state of future traffic. Third, the reflection on the iterative supervised learning approach and how this converges to an optimal policy. Fourth, the findings of naive perimeter control, perimeter control with an optimizer, and finally perimeter control with max pressure. The chapter concludes with a summary of the findings.

## 5.1 Model validation

Model validation is making sure that the simulation model resembles the real world. As it is assumed that the simulation model is accurate the validation is done by comparing the performance to a benchmark policy. A standard static traffic light policy is chosen as a benchmark. Iterative supervised learning can outperform the benchmark policy in terms of mean speed. However, the DRL algorithm is not able to outperform the benchmark. For the DRL algorithm, every attempted combination of state formulation, action, hyperparameters, and traffic resulted in performance just below the benchmark. As the DRL algorithm can choose the same policy as the benchmark policy there seems to be something wrong. Thorough testing is done but the exact cause is not found. It is expected that the problem is difficult to understand for the DRL algorithm which combined with a weak, delayed reward signal leads to poor performance.

## 5.2 Results neural network performance

This subsection covers the result of the neural network training. Here multiple state definitions are used and their performance is evaluated with the goal of finding the best state definition. For each unique state, the neural network architecture and hyperparameters are optimized for a fair comparison. This optimization however did not yield significantly different results. All state definitions performed best using almost identical parameters. One exception was state 5 which was better with an alpha of 0.01 instead of 0.001. Although the change in performance is insignificant. For ease of reading a recap of the states from Section 4.3 are given below:

| Feature | State 1 | State 2 | State 3 | State 4 | State 5 | State 6 |
|---|---|---|---|---|---|---|
| Total lane occupancy | ✓ | ✓ | ✓ | | | |
| Average lane occupancy | ✓ | | ✓ | ✓ | | |
| Total vehicles | ✓ | ✓ | ✓ | ✓ | | |
| Total vehicles inside perimeter | ✓ | ✓ | ✓ | ✓ | | ✓ |
| Total vehicles outside perimeter | ✓ | ✓ | ✓ | ✓ | | ✓ |
| Average speed | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| Total queue length | ✓ | ✓ | ✓ | | | |
| Average queue length | ✓ | | ✓ | ✓ | | |
| Queue outside perimeter | ✓ | ✓ | ✓ | ✓ | ✓ | |
| Queue inside perimeter | ✓ | ✓ | ✓ | ✓ | ✓ | |
| Total travel time | ✓ | ✓ | ✓ | | | |
| Average travel time | ✓ | | ✓ | ✓ | | |
| Inflowing vehicles | ✓ | ✓ | ✓ | ✓ | | |
| Outflowing vehicles | ✓ | ✓ | ✓ | ✓ | | |
| Queue length for every traffic light | | | ✓ | | | |

Table 7: Comparison of state formulations and their features.

State formulations 2 and 4 consist mostly of state 1 but the average values and total values were removed respectively if the alternative (total and average respectively) is present. State formulation 3 is an extension of state formulation 1 where the queue for each traffic light is added as a vector. States 5 and 6 are both parts small parts of state 1. After optimizing the neural network architecture and hyperparameters the following results were found. Note that the data is created using the same traffic.

Table 8: MAE, RMSE, and MAPE results for different state formulations.

| Model | MAE | RMSE | MAPE % |
|---|---|---|---|
| State 1 | 1.97 | 2.99 | 25.6 |
| State 2 | 1.84 | 2.66 | 22.4 |
| State 3 | 2.21 | 3.12 | 25.2 |
| State 4 | 1.79 | 2.62 | 21.1 |
| State 5 | 2.31 | 2.96 | 24.9 |
| State 6 | 2.08 | 3.09 | 25.2 |

Table 8 shows the mean absolute error (MAE), root mean square error (MRSE), and mean absolute percentage error (MAPE). The results show that state formulations 2 and 4 seem to be the best performing out of all variations tried. These two states are both similar to number one except for not having the combination of both "average of x" and "total of x" for the values that include both. It seems that having both these values in the state is redundant and decreases performance. It can also be observed that adding a lot of additional information such as the queue per traffic light decreases performance. Smaller state spaces with limited information such as states 5 and 6 give too little information. Because the performance of state space 4 is the best it has been chosen as the general state space for supervised learning and deep reinforcement learning for this problem. This

is the same state that is given in Section 3.2.



Pictures

Figure 9: Parity plot of the validation data using state 4.

Figure 9 shows how accurately the neural network can predict the labels. Most predicted values are close to the true values. However, there are a decent amount of outliers.



(a) Cumulative Error Distribution Plot using state 4.



(b) Density Plot of Prediction Errors using state 4.

Figure 10: Additional performance evaluation of the trained neural network using state 4.

Figure 10 shows how likely errors are expected to occur and the degree of the error as an additional measure of performance. Figure 10a shows that 70% of the values have an error margin below 20% and 90% of the data has an error margin below 40%.

The performance of the neural network might not seem great. This is in part due to the nature of the validation data, as the validation data is significantly different from the training data. The general relationship between the input state and the label seems to be understood by the neural network. When using similar traffic light settings and traffic intensities the parity plot has a much tighter fit and the error also reduces significantly.

## 5.3   Results iterative supervised learning

This section covers the results of applying a neural network to assign phase durations during the traffic simulation. After a lot of simulations are done the results can be found in Table 9. Here it can be seen that the error decreases over all performance parameters. This means that the iterative approach helps in improving the initial mediocre performance. This is done for both state 1 and state 4. This also shows that state formulation 4 remains superior to state formulation 1, even when fine-tuning is done to increase performance.

| Iteration | State 1 | | | State 4 | | |
|---|---|---|---|---|---|---|
| | MAE | RMSE | MAPE (%) | MAE | RMSE | MAPE (%) |
| 1 | 1.97 | 2.99 | 25.6 | 1.79 | 2.62 | 21.1 |
| 2 | 1.18 | 1.94 | 15.4 | 0.88 | 1.28 | 8.4 |
| 3 | 0.87 | 1.37 | 10.0 | 0.97 | 1.52 | 9.0 |
| 4 | 0.86 | 1.44 | 8.24 | 1.04 | 1.62 | 10.6 |
| 5 | 0.91 | 1.56 | 9.20 | 0.79 | 1.17 | 6.3 |

Table 9: MAE, RMSE, and MAPE values for state 1 & 4 per iteration of supervised learning.

In addition to the performance of predicting labels accurately, the performance of choosing good phase durations is also shown. Here the average speed over all data points is shown per iteration for both state 1 and state 4. This can be seen below in Figure 11. Note that iteration one is using a constant policy and is thus the same for both state 1 and state 4. Iteration two and onward use the model that was trained based on the data collected in the previous iteration.



Figure 11: Average speed per iteration of neural network updates for state 1 and state 4.

The results from Figure 11 show that the average speed over the simulation can be improved using the iterative supervised learning approach. Using the trained neural network to determine the policy yields a better average speed over many simulations. At the 5th iteration, the cumulative mean speed over five periods of the simulation is increased by more than 1 m/s with respect

to the starting random static policy. This is a percentage improvement of 5-8% which is not a lot. Compared to the standard static baseline policy the improvements are even smaller with the average speed being 0.5 m/s faster for the best-performing state. This is a percentage increase of 0-3% which is bad considering the base policy involves no intelligence. Noticeable in this graph is the decrease in performance for one iteration. This happens when a prediction of the neural network is significantly off and as a result, makes a really bad move. Consequently, this bad move causes severe congestion to the point that the overall traffic becomes stuck. In this case, it is no longer possible to fix the traffic. The result is an extremely low mean speed for that simulation run.

Except for the one iteration which decreases performance the other iterations all slowly increase performance. This performance is still on an upward trend at the 5th iteration, showing that more iterations and fine-tuning of the neural network can continue to improve the mean speed. More iterations have not been done due to the significant amount of time required to compute as creating the training and validation data takes a day per state.



(a) Cumulative Error Distribution Plot using state 4 at the 5th iteration.

(b) Parity plot of validation data using state 4 at the 5th iteration.

Figure 12: Performance of state 4 at the 5th iteration.

At the 5th iteration of fine-tuning for state 4 the error distribution has shifted as seen in Figure 12a. Now 85% of the data has an error of 10% or less. Originally this was 70% of the data having an error of 20% or less. Additionally, the parity curve is now much tighter around the identity line as seen in Figure 12b.

The difference between the performance of each state is also still prevalent. State 4 consistently outperforms state 1. However, despite state 1 being one of the worst states initially in terms of MAE, RMSE, and MAPE it still manages to improve over iterations. The errors are also reduced significantly as seen in Table 9. After fine-tuning the errors for both state 1 and state 4 have similar error values.

With the iterative supervised learning approach able to improve the mean cumulative speed over 5 cycles it should also be possible to apply DRL. These results will be discussed in the next section.

## 5.4 Results deep reinforcement learning

For the DRl model, a large combination of states, actions, and hyperparameters are used. The process to get this working took a lot of effort showing the complexity of correctly applying a DRL solution. The results seen in Figure 13 are the best results that are found using the naive DRL approach. Some learning can be observed as the reward increases compared to the starting point. However, this increase is not a lot and the mean speed of the best-performing iteration is still slightly below the standard static phase durations. It is assumed that the problem is too complicated to grasp. The state representation is very broad and the reward signal is delayed due to the nature of traffic. The combination of these factors is most likely the cause of the poor performance.

Changing from PPO to TD3 does not provide significant changes to the performance. Due to the similar results of PPO and TD3, it is assumed that the DRL algorithm itself is not the problem. What stands out in Figure 13 is that the line for PPO continues much longer. This is done to see if further improvements would be found if the simulation would be running for a very long time. This is not done for TD3 as it takes over a day to compute. The Figure also shows that both PPO and TD3 initially have a significant decrease in performance as they update their policy. This is a common occurance in DRL for multiple reasons such as exploration leading to poor results, value function instability where the learned value function is initially inacurate, and a poor formulation of the reward function which might encourage short term sub-optimal decisions. In this research the reason is most likely the initial inaccuracy in the value function. This is also shown in the iterative supervised learning approach that the initial prediction accuracy of a state-action pair is low. This low accuracy could lead to bad moves and thus a drop in reward. This change in reward is much smaller for PPO than for TD3, showing the effect of the clipping mechanism of PPO to prevent large sudden policy updates. Overall, the DRL learns over time but plateaus at a lower mean speed than the baseline policy which means that the performance is very poor.



Figure 13: Episode reward for TD3 and PPO using the naive action.

As the naive DRL approach shows such poor performance the focus of this research shifts more

towards improving the naive method as well as the supervised learning approaches instead of increasing the complexity of a solution method that is not working. Because of this the final solution method is not worked out and thus no results can be shown. The final solution method should in theory work well and fits all the criteria for the problem but unfortunately due to lack of time, it was not worked out.

## 5.5 Conclusion of the results

The neural network trained on regular traffic data performs fairly well on traffic states that are not shown before with 70% of the values having less than 20% error. The prediction error is significantly reduced by fine-tuning with 85% of the data having an error smaller than 10%. Iterative learning can be used to choose the phase durations and by doing so will improve the average speed by 5-8% compared to a static random policy or 0-3% compared to a standard static policy. This process takes a lot of time but the bottleneck is the simulation speed, not the neural network choosing the phase durations. The DRL approach is very challenging to implement correctly. After a lot of effort, the best results are still worse than the standard base traffic phase durations by a small margin.

# 6    Discussion

The results of this research show that it is possible to predict future congestion. However, the accuracy is limited with a MAPE ranging from 21-25% depending on the state formulation. The importance of providing a neural network with enough but also not too much information is shown by experimenting with different state formulations. After additional fine-tuning of the original neural network, the performance does increase significantly with a newly found MAPE between 8-10%. Now that the error is reduced significantly the neural network can be applied in the real world. As the computational bottleneck is collecting data and not generating outputs neural networks can be used in real time. This could be an interesting application for Technolution to provide.

It would be interesting to perform more iterations of updating the neural network to see how much it can still be improved. The error metrics seem to not improve much anymore but the mean cumulative speed is still increasing per iteration. It is assumed that eventually, the improvements will plateau but it is unknown whether that is after one or two more iterations or only after another ten or more iterations. This takes a lot of computational time but depending on the result it might be worthwhile if it offers a significant increase in the mean speed.

The iterative supervised learning approach also struggles to significantly outperform the fixed phase duration approach despite being able to predict future traffic relatively well. As the neural network can predict traffic, the problem might be in understanding the effects of the action. However, changing the action did not improve performance. Perhaps the problem is not the formulation of the action but rather the delayed reward signal as it takes a while for a traffic policy to have a significant effect on a network. This does showcase the challenges of working with neural networks, it is difficult to understand why they fail.

In contrast to the positive results achieved in predicting traffic states, the performance of the DRL model falls short of expectations. It consistently underperforms when compared to a standard scenario with fixed phase durations. Despite extensive testing and adjustments, pinpointing the exact reasons for this underperformance remains challenging. Potential causes may include issues with the reward structure or state representation, although numerous iterations and refinements have been attempted leading to the belief that the problem may be too complex for the DRL algorithm to learn effectively. This is further evidenced by the policy and value losses observed in TensorBoard, which remain exceedingly high and fail to decrease over time, which means that the algorithm struggles to accurately interpret the environment and understand the impact of its actions.

A possible way to improve the DRL performance would be to try different reward structures in addition to the ones tried before. Kampitakis & Vlahogianni (2024) for example used a scaled combination of trip completion rate, mean speed, and punishment for each simulation step taken to push the algorithm to empty the network as fast as possible as they had a finite amount of traffic. The trip completion rate is a great performance metric but is impossible to determine in reality. This is also an issue for the state representation, determining all values of the state accurately in reality is quite a challenge. Many traffic detection systems could be placed in the city, which will be costly. Even then, measuring for example how many cars are in the city as well as the amount of cars approaching the city is difficult. Similarly, determining the mean travel time for cars is nearly

impossible.

For the DRL approach, the focus is on policy-based DRL where the algorithm aims to find the policy that maximizes the reward directly. Alternatively, a value-based DRL approach could be used, which instead estimates a value function to predict the reward for a given state-action pair. Traditional value-based DRL methods such as Q-learning would not be suitable given the dimension of the action space. However, Deep Q-learning (DQN) could be effective in using the strengths of neural networks to learn the value function. Although, to use DQN the problem would have to be reformulated to include a discrete action space limiting the possible traffic light settings. The approach might also struggle with the complex state space of this problem which was already seen in the policy-based approach implemented in this research. Nonetheless, it could be interesting to try given that the value-based DRL approach is more sample efficient than the policy-based approach meaning that it could require less computational time which is a big issue as will be discussed in the next section.

One of the challenges encountered when trying to improve performance is that it takes a long time to see the result of a change. This is a significant limiting factor in exploring possible solutions to the performance issue. A large network that represents a crowded city is simply very slow to simulate. The network used in this research is not even that large containing only 25 intersections. If a large city would have to be simulated the computational times would be extremely long. A possible way to reduce the computational problem would be pre-training to reduce the iterations required for the DRL to reach convergence.

# 7 Conclusion

This section summarizes the research and finishes with an answer to the main research question. The literature highlights the complexity and stochastic nature of traffic, making it challenging to optimize and leading to time-intensive training for preventative solution methods. Reinforcement learning has emerged as a promising approach in traffic management due to its ability to make sequential decisions under uncertainty. However, most applications focus on single traffic light optimization, which, while effective, is not suitable for decision support given effective decision support requires infrequent actions and interpretable solutions.

To provide decision support an alternative to micro-level optimization is needed, in this case, global optimization over a larger time scale. A suitable method for this is perimeter control, which modifies the inflow and outflow of a city to reduce congestion. This technique requires infrequent actions and is intuitive to understand, making it excellent for decision support. The problem with only using perimeter control is that it modifies all traffic lights equally, which only works if traffic is uniform. Since in reality traffic is not uniform, an additional control layer is required, being max pressure. Max pressure ensures that green time is not wasted at traffic lights, thus providing local optimization. Deep reinforcement learning can be used to make the decision for perimeter control, determining the in and out flow modifications. In the mean time, a solver solves a set of linear equations to implement max pressure. This combined approach of global and local optimization should be able to provide decision support to help traffic engineers reduce congestion.

To work towards the final DRL model a few intermediate steps are taken to work towards the final DRL model: First, train a neural network to determine the effectiveness in predicting traffic congestion. Second, using the trained neural network to make decisions, to see if the trained neural network can not only learn the state but also find a good action in a given state. Lastly, creating a DRL model that implements perimeter control.

The results show that neural networks can be used for predicting traffic congestion, but unfortunately, not yet very accurate. Without any fine-tuning the neural network can predict traffic with a MAPE of 21-25%. With such a high margin of error at this point, it is not yet very practical. Alternatively, when using the neural network to make decisions, by adjusting the phase durations, the results show that the neural network can improve the mean speed. After multiple iterations, the average speed is improved by 5-8% with respect to the original random static policy and 0-3% with respect to the standard static policy. It is a computationally intense approach as the neural network improves per iteration, which takes a long time to simulate. Nevertheless, the newfound data can be used to fine-tune the original neural network, which improves the accuracy significantly resulting in a MAPE of 6-8% after fine-tuning the original neural network. This fine-tuned neural network can be effectively used as the results are fas and now also accurate.

The results of the DRL approach are not promising. The performance remains lower than the static base traffic light configuration despite experimenting with different rewards, state formulations, DRL algorithms, and hyperparameters. This approach is also computationally intense since many episodes are required, which takes a long time to simulate. In theory, the DRL approach should work well but sadly did not in this research. The problem assumed by the author is likely due to the DRL algorithm having difficulty understanding the effects of the actions on the state. A

reason for the author's assumption is the consistently high-value loss and policy loss.

To reflect on the research question "**How suitable is deep reinforcement learning for providing real-time decision support to traffic engineers for reducing congestion across multiple intersections in an urban traffic network?**", DRL is in theory very promising for sequential decision-making under uncertainty and has been shown to improve traffic in various cases. However, this research experienced a lot of difficulty with successfully applying DRL to optimize traffic. The performance is worse than a standard static policy making it unfit to be applied in the real world unless more time is invested. In addition, training the DRL algorithm itself is computationally intense. Although once trained, the DRL algorithm can return actions quickly. An alternative solution found to provide decision support is using neural networks to predict traffic. After fine-tuning the neural network can predict future congestion with good accuracy. Traffic engineers can input the state of traffic and the neural network can then return a prediction on what would happen in the future, giving valuable insights to traffic engineers.

# 8 Recommendations and future research

For future work, it would be recommended to use methods to reduce computational times as this was one of the main challenges. For running the simulation the following actions could be done to reduce computational time: Use smaller networks, less traffic, additional code optimization, or a faster traffic simulation software. To reduce the iterations required for training a DRL algorithm pre-training can be done. Alternatively, using advanced architectures like SimBa (Simplicity Bias) can be applied to enhance sample efficiency and reduce the episodes required for convergence (Lee et al. 2024). Kampitakis & Vlahogianni (2024) used pre-training based on behavior cloning to reduce the number of iterations required to reach convergence of the DRL algorithm. The idea is to learn a policy directly from a set of expert demonstrations. The goal is to teach the agent to replicate the expert's behavior. This can then be used as a strong starting policy. Another approach is to use better hardware that supports more parallel environments or GPU acceleration.

A difficult problem encountered during this research was understanding why the DRL decides on specific traffic light configurations. In the future, it could be possible to use a real-world network with traffic based on real data. Ideally, this network will already be optimized by a traffic expert. The DRL algorithm would then have to work towards the best solution which is assumed to be the one provided by the traffic expert. This could help understand the learning steps and could also provide insights to traffic experts on how artificial intelligence approaches this problem.

All in all, it is recommended that Technolution continues research into DRL as it is promising in theory. However, if they want to implement DRL they need to be aware that significant time and expertise is required. An alternative solution is to train neural networks to predict traffic congestion, which traffic engineers can use to gain insights. Although this is not the original goal it could be a useful tool that is much simpler to provide to traffic engineers.

## 9 References

## References

[1] Technolution Move. (2024). MobiMaestro - Technolution Move. Retrieved November 20, 2024, from `https://www.technolution.com/move/mobimaestro/`

[2] Inoue, D., Yamashita, H., Aihara, K., & Yoshida, H. (2024). Traffic signal optimization in large-scale urban road networks: an adaptive-predictive controller using Ising models. arXiv.org. `https://arxiv.org/abs/2406.03690`

[3] Li, Y. (2018) Deep Reinforcement Learning `https://arxiv.org/abs/1810.06339`

[4] Wyner, A. (2023). Sustainable safety: The dutch approach to safe road design `https://www.ippi.org.il/sustainable-safety-the-dutch-approach-to-safe-road-design/`

[5] Naeem, M., Rizvi, S. T. H., & Coronato, A. (2020). A Gentle Introduction to Reinforcement Learning and its Application in Different Fields. IEEE Access, 8, 209320–209344, `https://doi.org/10.1109/access.2020.3038605`

[6] Allied Market Research, (2023). Reinforcement Learning Market Size, share, Competitive landscape and Trend analysis Report by deployment mode, by enterprise size, by end user: Global Opportunity Analysis and Industry Forecast, 2023-2032. Allied Market Research. `https://www.alliedmarketresearch.com/reinforcement-learning-market-A229407#:~:text=Reinforcement%20Learning%20Market%20Statistics%3A%202032,41.5%25%20from%202023%20to%202032`

[7] Wu, C., Kreidieh, A., Parvate, K., Vinitsky, E., & Bayen, A. M. (2017). Flow: Architecture and benchmarking for reinforcement learning in traffic Control. arXiv (Cornell University). `https://arxiv.org/pdf/1710.05465.pdf`

[8] Zhang, H., Feng, S., Liu, C., Ding, Y., Zhu, Y., Zhou, Z., Zhang, W., Yu, Y., Jin, H., & Li, Z. (2019). CityFlow: A Multi-Agent Reinforcement learning environment for large scale city traffic scenario. In Proceedings of the 2019 World Wide Web Conference (WWW '19), May 13–17, 2019, San Francisco, CA, USA. `https://doi.org/10.1145/3308558.3314139`

[9] Qiang, W., & Zhongli, Z. (2011). Reinforcement Learning Model,Algorithms and Its Application. International Conference on Mechatronic Science, Electric Engineering and Computer, Jilin, China, pp. 1143-1146, `https://doi.org/10.1109/MEC.2011.6025669`

[10] Creswell, J. W. (2023). Research design: Qualitative, Quantitative and Mixed Methods approaches. Los Angeles: SAGE Publications Inc. From `http://www.revistacomunicacion.org/pdf/n3/resenas/research_design_qualitative_quantitative_and_mixed_methods_approaches.pdf`

[11] McCluskey, T., Kotsialos, A., Mller, J. P., Klgl, F., Rana, O. F., & Schumann, R. (2016). Autonomic Road transport support systems. In Springer eBooks. `https://doi.org/10.1007/978-3-319-25808-9`

[12] Varaiya, P. (2013). Max pressure control of a network of signalized intersections. Transportation Research Part C Emerging Technologies, 36, 177–195. `https://doi.org/10.1016/j.trc.2013.08.014`

[13] Technolution Move. (2023). AFM - Technolution Move. `https://www.technolution.com/move/cases/afm/`

[14] Rasheed, F., Yau, K. A., Noor, R. M., Wu, C., & Low, Y. C. (2020). Deep Reinforcement Learning for Traffic Signal Control: a review. IEEE Access, 8, 208016–208044. `https://doi.org/10.1109/access.2020.3034141`

[15] Abdulhai, B., Pringle, R., & Karakoulas, G. J. (2003). Reinforcement learning for true adaptive traffic signal control. Journal of Transportation Engineering, 129(3), 278–285. `https://doi.org/10.1061/(asce)0733-947x(2003)129:3(278`

[16] Genders, W., & Razavi, S. (2016). Using a deep reinforcement learning agent for traffic signal control. arXiv (Cornell University). `http://export.arxiv.org/pdf/1611.01142`

[17] Zheng, G., Zang, X., Xu, N., Wei, H., Yu, Z., Gayah, V. V., Xu, K., & Li, Z. (2019). Diagnosing reinforcement learning for traffic signal control. arXiv (Cornell University). `https://arxiv.org/pdf/1905.04716.pdf`

[18] Li, Z., Yu, H., Zhang, G., Dong, S., & Xu, C. (2021). Network-wide traffic signal control optimization using a multi-agent deep reinforcement learning. Transportation Research Part C Emerging Technologies, 125, 103059. `https://doi.org/10.1016/j.trc.2021.103059`

[19] Ouallane, A. A., Bahnasse, A., Bakali, A., & Talea, M. (2022). Overview of Road Traffic Management Solutions based on IoT and AI. Procedia Computer Science, 198, 518–523. `https://doi.org/10.1016/j.procs.2021.12.279`

[20] Ravish, R., & Swamy, S. R. (2021). Intelligent Traffic Management: A review of challenges, solutions, and future perspectives. Transport and Telecommunication, 22(2), 163–182. `https://doi.org/10.2478/ttj-2021-0013`

[21] Cong, Z., De Schutter, B., & Babuška, R. (2013) Ant colony routing algorithm for freeway networks. Transportation Research Part C: Emerging Technologies, 37, 1-19. `https://www.sciencedirect.com/science/article/pii/S0968090X13001885`

[22] Shi, Y., Qi, Y., Lv, L., & Liang, D. (2021). A Particle Swarm Optimisation with Linearly Decreasing Weight for Real-Time Traffic Signal Control. Machines, 9(11), 280. `https://doi.org/10.3390/machines9110280`

[23] Tsitsokas, D., Kouvelas, A., & Geroliminis, N. (2023). Two-layer adaptive signal control framework for large-scale dynamically-congested networks: Combining efficient Max Pressure with Perimeter Control. Transportation Research. Part C, Emerging Technologies, 152, 104128. `https://doi.org/10.1016/j.trc.2023.104128`

[24] Tsitsokas, D., Kouvelas, A., & Geroliminis, N. (2021). Modeling and optimization of dedicated bus lanes space allocation in large networks with dynamic congestion. Transportation Research Part C Emerging Technologies, 127, 103082. `https://doi.org/10.1016/j.trc.2021.103082`

[25] Tampubolon, H., & Hsiung, P. (2018). Supervised Deep Learning Based for Traffic Flow Prediction. 2018 International Conference on Smart Green Technology in Electrical and Information Systems (ICSGTEIS), Bali, Indonesia, pp. 95-100, doi: `10.1109/ICSGTEIS.2018.8709102`

[26] Sanagapati, P. (2019). What is Dropout Regularization? Find out:). Kaggle. Retrieved November 12, 2024, from `https://www.kaggle.com/code/pavansanagapati/what-is-dropout-regul arization-find-out`

[27] Patil , P. (2022). Applications of Deep Learning in Traffic Management: A Review. International Journal of Business Intelligence and Big Data Analytics, 5(1), 16–23. Retrieved November 12, 2024, from `https://research.tensorgate.org/index.php/IJBIBDA/article/view/26`

[28] Wei, H., Chen, C., Wu, K., Xu, K., Gayah, V., & Li, Z. (2019). Presslight: Learning max pressure control to coordinate traffic signals in arterial network. In Proceedings of the 25th ACMSIGKDD International Conference on Knowledge Discovery & Data Mining. `https://ojs.aaai.org/index.php/AAAI/article/view/5744/5600`

[29] Hegde, S., Premasudha, B., Hooli, A., Akshay, M. (2024). A Review on Smart Traffic Management with Reinforcement Learning. In: Yang, XS., Sherratt, S., Dey, N., Joshi, A. (eds) Proceedings of Ninth International Congress on Information and Communication Technology. ICICT 2024 2024. Lecture Notes in Networks and Systems, vol 1004. Springer, Singapore. https://doi.org/10.1007/978-981-97-3305-7$_3$7

[30] Cong, Z., De Schutter, B., & Babuška, R. (2013). Ant Colony Routing algorithm for freeway networks. Transportation Research Part C Emerging Technologies, 37, 1–19. `https://doi.org/ 10.1016/j.trc.2013.09.008`

[31] Han, Y., Wang, M., & Leclercq, L. (2023). Leveraging reinforcement learning for dynamic traffic control: A survey and challenges for field implementation. Communications in Transportation Research, 3, 100104. `https://doi.org/10.1016/j.commtr.2023.100104`

[32] Lin, Y., Dai, X., Li, L., & Wang, F. (2018). An efficient deep reinforcement learning model for urban traffic control. arXiv (Cornell University). `http://export.arxiv.org/pdf/1808.01876`

[33] Han, Y., Wang, M., & Leclercq, L. (2023). Leveraging reinforcement learning for dynamic traffic control: A survey and challenges for field implementation. Communications in Transportation Research, 3, 100104. `https://doi.org/10.1016/j.commtr.2023.100104`

[34] Li, Z., Xu, C., Zhang, G. (2021). A Deep Reinforcement Learning Approach for Traffic Signal Control Optimization `https://arxiv.org/abs/2107.06115`

[35] , L., Babuska, R., & De Schutter, B. (2010). Multi-agent reinforcement learning: An overview. In D. Srinivasan, & L. C. Jain (Eds.), Innovations in multi-agent systems and applications - 1 (pp. 183-221). Springer.

[36] Li, Y. (2018).Deep Reinforcement Learning. urlhttps://arxiv.org/abs/1810.06339

[37] Wang, X., Wang, S., Liang, X., Zhao, D., Huang, J., Xu, X., Dai, B., & Miao, Q. (2024). Deep Reinforcement Learning: a survey. IEEE Transactions on Neural Networks and Learning Systems, 1–15. `https://doi.org/10.1109/tnnls.2022.3207346`

[38] Walraven, E., Spaan, M. T. J., & Bakker, B. (2016). Traffic flow optimization: A reinforcement learning approach. Engineering Applications of Artificial Intelligence, 52, 203–212. `https://doi.org/10.1016/j.engappai.2016.01.001`

[39] Paul, A., Mitra, S. (2020). Deep Reinforcement Learning based Traffic Signal optimization for Multiple Intersections in ITS, IEEE International Conference on Advanced Networks and Telecommunications Systems (ANTS), New Delhi, India, 2020, pp. 1-6. `https://doi.org/10.1109/ANTS50601.2020.9342819`

[40] Gu, J., Lee, M., Jun, C., Han, Y., Kim, Y., Kim, J. (2021). Traffic Signal Optimization for Multiple Intersections Based on Reinforcement Learning. Appl. Sci. 2021, 11, 10688. `https://doi.org/10.3390/app112210688`

[41] Muresan, M., Pan, G., & Fu, L. (2021). Multi-Intersection Control with Deep Reinforcement Learning and Ring-and-Barrier Controllers. Transportation Research Record, 2675(4), 308-319. `https://doi.org/10.1177/0361198120980321`

[42] Zhu, R., Ding, W., Wu, S., Li, L., Lv, P., & Xu, M. (2023). Auto-learning communication reinforcement learning for multi-intersection traffic light control. Knowledge-based Systems, 275, 110696. `https://doi.org/10.1016/j.knosys.2023.110696`

[43] Wang, Y., Xu, T., Niu, X., Tan, C., Chen, E., & Xiong, H. (2022). STMARL: A Spatio-Temporal Multi-Agent Reinforcement Learning Approach for Cooperative Traffic Light Control. IEEE Transactions on Mobile Computing, 21(6), 2228–2242. `https://doi.org/10.1109/tmc.2020.3033782`

[44] Tsitsokas, D., Kouvelas, A., & Geroliminis, N. (2022). Efficient Max-Pressure traffic management for Large-Scale congested urban networks. Infoscience. `https://infoscience.epfl.ch/record/301369`

[45] Handboek verkeerslichtenregelingen 2022. (2022). Kennisbank.crow.nl. Retrieved June 4, 2024, from `https://kennisbank.crow.nl/public/VERMGT/Handboek_Verkeerslichtenregelingen_2022/Minimale_groentijd_en_groenknippertijd_voor_voetgangers/114372`

[46] Yadav, N., Mathur, P., Kumar, P., & Prasad, N. (2022). Road Traffic Signal Design By Using Webster Method. International Journal of Scientific Engineering and Applied Science (IJSEAS) – Volume-8, Issue-2. `https://ijseas.com/volume8/v8i2/IJSEAS202202102.pdf`

[47] Van Heeswijk, W., (2023). Proximal Policy Optimization (PPO) explained - towards data science. Medium. `https://towardsdatascience.com/proximal-policy-optimization-ppo-explained-abed1952457b`

[48] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal Policy Optimization Algorithms. `https://arxiv.org/abs/1707.06347`

[49] Raffin, A., Hill, A., Ernestus, M., Gleave, A., Kanervisto, A., & Dormann, N. (2019). Stable-Baselines3: Reliable Reinforcement Learning Implementations. GitHub repository. `https://github.com/DLR-RM/stable-baselines3`

[50] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 12, 2825–2830. `http://scikit-learn.org/stable`

[51] Kampitakis, E., & Vlahogianni, E. (2024). Unified cooperative management system: a framework to integrate the control strategies. Manuscript submitted for publication.

[52] Lee, H., Hwang, D., Kim, D., Kim, H., Tai, J., Subramanian, K., Wurman, P., Choo, J., Stone, P., Seno, T. (2024). SimBa: Simplicity Bias for Scaling Up Parameters in Deep Reinforcement Learning. `https://arxiv.org/abs/2410.09754`

# A Appendix

## A.1 Additional figures



Figure A.1: Parity plot of the validation data using state 1.



(a) Cumulative Error Distribution Plot using state 1.

(b) Density Plot of Prediction Errors using state 1.

Figure A.2: Additional performance evaluation of the trained neural network using state 1.

Figure A.3: Parity plot of the validation data using state 2.



(a) Cumulative Error Distribution Plot using state 2.

(b) Density Plot of Prediction Errors using state 2.

Figure A.4: Additional performance evaluation of the trained neural network using state 2.

Figure A.5: Parity plot of the validation data using state 3.



(a) Cumulative Error Distribution Plot using state 3.

(b) Density Plot of Prediction Errors using state 3.

Figure A.6: Additional performance evaluation of the trained neural network using state 3.

Figure A.7: Parity plot of the validation data using state 5.



(a) Cumulative Error Distribution Plot using state 5.



(b) Density Plot of Prediction Errors using state 5.

Figure A.8: Additional performance evaluation of the trained neural network using state 5.

Figure A.9: Parity plot of the validation data using state 6.



(a) Cumulative Error Distribution Plot using state 6.



(b) Density Plot of Prediction Errors using state 6.

Figure A.10: Additional performance evaluation of the trained neural network using state 6.

## A.2 Python code

```python
1  import traci
2  import sumolib
3  # this generates the actual trips used in the warm up simulation
4  from generate_trips_NN import run_full_pipeline
5
6  """"This file will create the warm up state for the city simulations.
7      This will create a .xml file which can be loaded in as starting point for other simulati
8
9  def save_traffic_state(sumo_cmd, state_file, pre_sim_duration=600):
10     """
11     Runs a pre-simulation phase and saves the traffic state to a file.
12
13     Parameters:
14     - sumo_cmd: SUMO command for starting the simulation.
15     - state_file: Path to save the state file.
16     - pre_sim_duration: Duration of the pre-simulation phase (in seconds).
17     """
18     traci.start(sumo_cmd)
19     for step in range(pre_sim_duration):
20         traci.simulationStep()
21     traci.simulation.saveState(state_file)
22     traci.close()
23     print(f"Successfully saved the state to {state_file}")
24
25
26     input_file = state_file
27
28     # Read the original file, filtering out lines that start with '<flowState id="'
29     with open(input_file, "r") as file:
30         lines = file.readlines()
31
32     # Write back only lines that do not start with '<flowState id="'
33     with open(input_file, "w") as file:
34         for line in lines:
35             if not line.strip().startswith('<flowState id="'):
36                 file.write(line)
37
38  if __name__ == "__main__":
39      network_file = "sumo_rl/nets/PC-Alex/city_network.net.xml"
40      output_file = "sumo_rl/nets/PC-Alex/warm_up_routing.rou.xml"
41      fringe_factor = 1  # Adjust fringe factor as needed
42      begin_time = 0
43      end_time = 600
44      trip_count = 2200
45
46      # step 1: generate routing
47      # Run the full pipeline with the given variables to create traffic
```

```
48        run_full_pipeline(network_file, output_file, fringe_factor, begin_time, end_time, trip_c

49

50        # Step 2: Save the initial traffic state
51        sumoBinary = sumolib.checkBinary('sumo-gui')
52        sumo_cmd = [
53        sumoBinary,
54        "-n", "sumo_rl/nets/PC-Alex/City_network.net.xml",
55        "-r", output_file ,
56        "--fcd-output", "fcd.xml",
57        "--vehroute-output", "vehroutes.xml",
58        "--default.speeddev", "0",
59        "--no-step-log",
60        "--summary-output", "summary_output.xml"]
61        state_file = "City_warm_up.xml"
62        save_traffic_state(sumo_cmd, state_file, end_time)
```

```
1    import subprocess
2    import os
3    import xml.etree.ElementTree as ET
4    import math
5    import random
6    import numpy as np
7    import time

8

9    " this file will generate trips for the simulation creating a route file"

10

11   def generate_center_focused_trips(network_file, output_file, center_nodes, outside_nodes_sta
12        """
13        Generates random trips focusing more on nodes near the center of the network.

14

15        Parameters:
16        - network_file: Path to the SUMO network file (.net.xml).
17        - output_file: Path to save the generated routes file (.rou.xml).
18        - center_nodes: List of nodes or edges near the center of the network to favor for start
19        - fringe_factor: The fringe factor to influence trip start and end locations.
20        - begin_time: The start time for generating trips.
21        - end_time: The end time for generating trips.
22        - trip_count: The number of trips to generate.
23        """
24        command = [
25            "python",
26            os.path.join(os.environ["SUMO_HOME"], "tools", "randomTrips.py"),
27            "-n", network_file,
28            "-r", output_file,
29            "--fringe-factor", str(fringe_factor),
30            "-b", str(begin_time),
31            "-e", str(end_time),
32            "-p", str((end_time - begin_time) / trip_count),
33            "--trip-attributes", "",
```

```
34              "--validate",
35          ]
36
37          # Generate trips with more focus on the center nodes
38          try:
39              # Create a temporary trips file with all trips and adjust start/end points
40              temp_trip_file = output_file.replace('.xml', '_temp.xml')
41              command.extend(["-o", temp_trip_file])
42
43              # Run the command to generate trips
44              subprocess.run(command, check=True)
45
46              # Parse the generated trips and adjust their start/end points
47              # this modifies the style and will only include a start and end point
48              adjust_trips_to_center(temp_trip_file, output_file, center_nodes, outside_nodes_star
49
50              print(f"Successfully generated center-focused trips and saved to {output_file}")
51          except subprocess.CalledProcessError as e:
52              print(f"Error generating trips: {e}")
53          except Exception as e:
54              print(f"An unexpected error occurred: {e}")
55
56  def adjust_trips_to_center(temp_trip_file, output_file, center_nodes, outside_nodes_starting
57      """
58      Adjusts trips in the temporary trip file to favor starting or ending near center nodes.
59
60      Parameters:
61      - temp_trip_file: The temporary trips file to adjust.
62      - output_file: Path to save the adjusted routes file (.rou.xml).
63      - center_nodes: List of nodes to favor for trip start/end points.
64      """
65
66      # Initialize dictionaries to track usage counts for each node
67      center_usage = {node: 0 for node in center_nodes}
68      starting_usage = {node: 0 for node in outside_nodes_starting}
69      exit_usage = {node: 0 for node in outside_nodes_exit}
70
71      def weighted_choice(node_usage):
72          """Select a node with a probability inversely proportional to its usage count, with
73          # Adjust the weight factor to prioritize nodes with lower usage more strongly
74          weight_factor = 3
75          total = sum((1 / (count + 1)) ** weight_factor for count in node_usage.values())
76          weights = [((1 / (count + 1)) ** weight_factor) / total for count in node_usage.valu
77          chosen_node = random.choices(list(node_usage.keys()), weights=weights, k=1)[0]
78          node_usage[chosen_node] += 1  # Update usage count for the chosen node
79          return chosen_node
80
81      tree = ET.parse(temp_trip_file)
82      root = tree.getroot()
83
```

```
84        # modify this to change the ratio of traffic going into the city
85        percentage_of_trips_into_city = 0.7
86        for trip in root.findall('trip'):
87            if random.random() < percentage_of_trips_into_city:
88                # Trip towards the center from the outside
89                to_node = weighted_choice(center_usage)
90                from_node = weighted_choice(starting_usage)
91            else:
92                # Trip from the center to the outside
93                from_node = weighted_choice(center_usage)
94                to_node = weighted_choice(exit_usage)
95
96            # Set the 'from' and 'to' for each trip
97            trip.set('from', from_node)
98            trip.set('to', to_node)
99
100       # Write adjusted trips to the final output file
101       tree.write(output_file)
102       os.remove(temp_trip_file)  # Clean up temporary file
103
104
105   def run_full_pipeline(network_file, output_file, fringe_factor=1, begin_time=600, end_time=3
106       """
107       Runs the full pipeline: generates trips and extracts unique routes.
108       Arg: network_file: the .xml file which contains the network
109       arg: output_file: file to write the assigned trips to
110       arg: fringe_factor: degree of traffic spawning on edges compared to center
111       arg: begin_time: start time for traffic
112       arg: trip_count: how many trips will be created
113       arg: traffic_dist: the type of traffic that will be created
114       """
115       print(f"Generating trips with output file: {output_file}")
116       # all the edges that are not on or outside the perimeter
117       center_nodes = ["X2Y4.5L", "X2Y4.5R", "X3Y4.5L", "X3Y4.5R", "X4Y4.5L", "X4Y4.5R", "X1.5Y
118                       "X2Y3.5L", "X2Y3.5R", "X3Y3.5L", "X3Y3.5R", "X4Y3.5L", "X4Y3.5R", "X1.5Y
119                       "X2Y2.5L", "X2Y2.5R", "X3Y2.5L", "X3Y2.5R", "X4Y2.5L", "X4Y2.5R", "X1.5Y
120                       "X2Y1.5L", "X2Y1.5R", "X3Y1.5L", "X3Y1.5R", "X4Y1.5L", "X4Y1.5R",
121                       ]
122
123       # all outside edges that go into the perim (valid starting point)
124       outside_nodes_starting = ["X1Y5.5L", "X2Y5.5L", "X3Y5.5L", "X4Y5.5L", "X5Y5.5L",
125                       "X0.5Y5D", "X0.5Y4D", "X0.5Y3D", "X0.5Y2D", "X0.5Y1D",
126                       "X1Y0.5R", "X2Y0.5R", "X3Y0.5R", "X4Y0.5R", "X5Y0.5R",
127                       "X5.5Y5U", "X5.5Y4U", "X5.5Y3U", "X5.5Y2U", "X5.5Y1U",
128                       ]
129
130       outside_nodes_exit = ["X1Y5.5R", "X2Y5.5R", "X3Y5.5R", "X4Y5.5R", "X5Y5.5R",
131                       "X0.5Y5U", "X0.5Y4U", "X0.5Y3U", "X0.5Y2U", "X0.5Y1U",
132                       "X1Y0.5L", "X2Y0.5L", "X3Y0.5L", "X4Y0.5L", "X5Y0.5L",
133                       "X5.5Y5D", "X5.5Y4D", "X5.5Y3D", "X5.5Y2D", "X5.5Y1D",
```

65

```
134                        ]
135
136      # Step 1: Generate random trips
137      # in this function adjust_trips_to_center is called wich modifies the trip file to have
138      generate_center_focused_trips(network_file, output_file, center_nodes, outside_nodes_sta
139
140      # step 2: The new center focussed trips include a start and end point but require a rout
141      generate_complete_routes(output_file, end_time, traffic_dist)
142
143      # Step 3: Extract unique routes from the generated trip file
144      input_route_file = output_file # input is the output file of generate_center_focused_tri
145      output_route_file = output_file   #only adjust the content, not the name
146
147
148  def generate_complete_routes(input_route_file, end_time=3600, traffic_dist='exp(0.0043)'):
149      """
150      Uses duarouter to generate complete routes with intermediate edges and writes them to a
151
152      Parameters:
153      - input_route_file: Path to the input trip file (.xml).
154      - output_route_file: Path to save the complete routes file (.rou.xml).
155      - temp_trip_file: Path for a temporary .trips.xml file for duarouter processing.
156      - end_time: End time for the flow generation.
157      - traffic_dist: Defines the traffic distribution type (poisson, uniform, etc.)
158      """
159
160      # Step 1: Create a temporary .trips.xml file with <trip> elements
161
162      # temp trip file to save the trips to
163      temp_trip_file = input_route_file.replace('.xml', '_temp_input.xml')
164      # this file will then be renamed to the input_route_file as this is the desired output f
165      # the original input_route_file will be deleted after it's information is used
166      temp_output_route_file = input_route_file.replace('.xml', '_temp_output.xml')
167
168      tree = ET.parse(input_route_file)
169      root = tree.getroot()
170
171      # Write temporary trips file
172      with open(temp_trip_file, 'w') as f:
173          f.write('<?xml version="1.0" encoding="UTF-8"?>\n<trips>\n')
174          for trip in root.findall('trip'):
175              trip_id = trip.get('id')
176              depart = trip.get('depart')
177              from_edge = trip.get('from')
178              to_edge = trip.get('to')
179              f.write(f'    <trip id="{trip_id}" depart="{depart}" from="{from_edge}" to="{to_
180          f.write('</trips>')
181
182      # Step 2: Run duarouter to compute routes with all intermediate edges
183      try:
```

```
184        subprocess.run([
185            "duarouter",
186            "-n", "sumo_rl/nets/PC-Alex/city_network.net.xml",
187            "-t", temp_trip_file,
188            "-o", temp_output_route_file,
189            "--ignore-errors",  # This prevents duarouter from stopping due to minor connect
190        ], check=True)
191        print(f"Routes successfully generated in {temp_output_route_file}")
192    except subprocess.CalledProcessError as e:
193        print(f"Error generating routes with duarouter: {e}")
194        return
195
196    # Step 3: Open the generated .rou.xml file and add flows if necessary
197    tree = ET.parse(temp_output_route_file)
198    root = tree.getroot()
199    with open(temp_output_route_file, 'w') as f:
200        f.write('<?xml version="1.0" encoding="UTF-8"?>\n<routes>\n')
201
202        # Define vehicle types outside of the routes
203        f.write('<vType id="CarA" accel="2.0" decel="4.5" sigma="0.5" length="5.0" maxSpeed=
204
205        begin_time = 600 # end time of the warm up period
206        # Write each unique route and associated flows
207        for i, vehicle in enumerate(root.findall('vehicle')):
208            route_id = vehicle.get('id')
209            edges = vehicle.find('route').get('edges')
210            f.write(f'    <route id="{route_id}" edges="{edges}"/>\n')
211
212            # Adding flow with custom period based on traffic_dist
213            # if period is equal to x it means there are x seconds between each vehicle creat
214            if traffic_dist == 'step':
215                # i loops over the vehicles, as more vehicles are created, the period betwee
216                if 0 <= i < 1000:
217                    period = 5 #light traffic
218                elif 1000 <= i < 2000:
219                    period = 3 # moderate traffic
220                else:
221                    period = 1  # heavy traffic
222            elif traffic_dist == 'sinusoidal':
223                # this creates a period with base x +/- 2 as math.sin(i / 1000.0) oscilliate
224                period = 5 + 3 * math.sin(i / 10)
225            elif traffic_dist.startswith("poisson"):
226                # Extract the values inside the parentheses
227                # slices up the input "poisson(x)", and takes the value between poisson( and
228                avg_rate = float(traffic_dist[len("poisson("):-1])
229                period = np.random.exponential(1 / avg_rate)
230            else:
231                period = traffic_dist
232
233            if i % 50 ==0:
```

```
234                     begin_time += 200
235
236                 # after a while the simulation starts getting too full, this reduces the input
237                 if begin_time > 1000:
238                     avg_rate = float(traffic_dist[len("exp("):-1])
239                     avg_rate = 0.85*avg_rate
240                     period = f"exp({avg_rate})"
241                 if begin_time > 1500:
242                     avg_rate = float(traffic_dist[len("exp("):-1])
243                     avg_rate = 0.75*avg_rate
244                     period = f"exp({avg_rate})"
245                 if begin_time > 2000:
246                     avg_rate = float(traffic_dist[len("exp("):-1])
247                     avg_rate = 0.7*avg_rate
248                     period = f"exp({avg_rate})"
249
250                 f.write(f'    <flow id="flow{i+1}" route="{route_id}" type="CarA" beg="{begin_ti
251             f.write('</routes>')
252
253
254
255
256         # Clean up the temporary trip file
257         os.remove(temp_trip_file)
258
259         def safe_remove(file_path, retries=5, delay=2):
260             """Attempts to remove a file with retries in case of permission issues."""
261             for attempt in range(retries):
262                 try:
263                     os.remove(file_path)
264                     print(f"Successfully removed {file_path}")
265                     break
266                 except PermissionError:
267                     print(f"PermissionError on attempt {attempt + 1} to remove {file_path}. Retr
268                     time.sleep(delay)
269             else:
270                 print(f"Failed to remove {file_path} after {retries} attempts.")
271                 raise PermissionError(f"Could not remove {file_path} after multiple attempts")
272         safe_remove(input_route_file)
273
274         time.sleep(0.1)
275         # Rename `output_route_file` to `output_file`
276         os.rename(temp_output_route_file, input_route_file)
277
278 if __name__ == "__main__":
279     network_file = "sumo_rl/nets/PC-Alex/city_network.net.xml"
280     output_file = "sumo_rl/nets/PC-Alex/unique_routes_NN.trips.xml"
281     fringe_factor = 1
282     begin_time = 600
283     end_time = 3600
```

68

```
284        trip_count = 4000
285
286        # Run the full pipeline with the given variables
287        run_full_pipeline(network_file, output_file, fringe_factor, begin_time, end_time, trip_c
288        #generate_center_focused_trips(network_file, output_file, center_nodes, fringe_factor, b
289
290        # Usage
291        input_route_file = "sumo_rl/nets/PC-Alex/unique_routes_NN.trips.xml"  # Path to your ini
292        output_route_file = "sumo_rl/nets/PC-Alex/unique_routes_NN.rou.xml"  # Path to save the
293        temp_trip_file = "sumo_rl/nets/PC-Alex/temp.trips.xml"  # Temporary trip file for duarou
294
295        generate_complete_routes(input_route_file, output_route_file, temp_trip_file)
```

```
1  import traci
2  import sumolib
3  # this generates the actual trips used in the war up simulation
4  from generate_trips_NN import run_full_pipeline
5
6  """"This file will create the warm up state for the city simulations.
7      This will create a .xml file which can be loaded in as starting point for other simulati
8
9  def save_traffic_state(sumo_cmd, state_file, pre_sim_duration=600):
10     """
11     Runs a pre-simulation phase and saves the traffic state to a file.
12
13     Parameters:
14     - sumo_cmd: SUMO command for starting the simulation.
15     - state_file: Path to save the state file.
16     - pre_sim_duration: Duration of the pre-simulation phase (in seconds).
17     """
18     traci.start(sumo_cmd)
19     for step in range(pre_sim_duration):
20         traci.simulationStep()
21     traci.simulation.saveState(state_file)
22     traci.close()
23     print(f"Successfully saved the state to {state_file}")
24
25
26     input_file = state_file
27
28     # Read the original file, filtering out lines that start with '<flowState id="'
29     with open(input_file, "r") as file:
30         lines = file.readlines()
31
32     # Write back only lines that do not start with '<flowState id="'
33     with open(input_file, "w") as file:
34         for line in lines:
35             if not line.strip().startswith('<flowState id="'):
36                 file.write(line)
```

69

```
37
38  if __name__ == "__main__":
39      network_file = "sumo_rl/nets/PC-Alex/city_network.net.xml"
40      output_file = "sumo_rl/nets/PC-Alex/warm_up_routing.rou.xml"
41      fringe_factor = 1  # Adjust fringe factor as needed
42      begin_time = 0
43      end_time = 600
44      trip_count = 2200
45
46      # step 1: generate routing
47      # Run the full pipeline with the given variables to create traffic
48      run_full_pipeline(network_file, output_file, fringe_factor, begin_time, end_time, trip_c
49
50      # Step 2: Save the initial traffic state
51      sumoBinary = sumolib.checkBinary('sumo-gui')
52      sumo_cmd = [
53      sumoBinary,
54      "-n", "sumo_rl/nets/PC-Alex/City_network.net.xml",
55      "-r", output_file ,
56      "--fcd-output", "fcd.xml",
57      "--vehroute-output", "vehroutes.xml",
58      "--default.speeddev", "0",
59      "--no-step-log",
60      "--summary-output", "summary_output.xml"]
61      state_file = "City_warm_up.xml"
62      save_traffic_state(sumo_cmd, state_file, end_time)
```

```
1   import traci
2   import libsumo
3   import numpy as np
4   import pandas as pd
5   import xml.etree.ElementTree as ET
6   import collections
7   import sumolib
8   from shapely.geometry import LineString, Point
9   import generate_trips_NN
10  import time
11  from joblib import load
12  import math
13
14  "file used to collect data when the NN is used for decision making (iterative supervised lea
15
16  class SumoDataCollector:
17      """Class to collect data for supervised learning"""
18
19      def __init__(self, sumo_cmd, use_libsumo=False):
20          self.sumo_cmd = sumo_cmd
21          self.data = []
22          self.avg_speeds_per_step = []
```

```python
        self.cycle_mean_speed = 0
        self.total_arrived_vehicles = 0
        self.perim_box = 0, -200, 800, 600 #modify this if the permiter changes
        self.min_green_time = 8
        self.controlled_tl_lights = []
        self.vehicle_ids = []
        self.traffic_light_data = {}
        # Import the correct module based on `use_libsumo var`
        if use_libsumo:
            import libsumo as sumo_interface
            print(' using libsumo')
        else:
            import traci as sumo_interface
            print(' using traci')
        self.sumo_interface = sumo_interface

        # Parse the .net.xml file
        tree = ET.parse('sumo_rl/nets/PC-Alex/City_network.net.xml')
        root = tree.getroot()

        # Get all lanes
        self.all_lanes = []
        for edge in root.findall('edge'):
            for lane in edge.findall('lane'):
                self.all_lanes.append(lane.get('id'))

    def create_traffic(self, traffic_dist, trip_count):
        """
        Generates the traffic flows for the simulation. This is called so that the NN can ha
        param: traffic_dist: intensity of traffic on the created routes as poisson arrival p
        param: trip_count: how many trips are present during the simulation (how busy the ne
        """
        network_file = "sumo_rl/nets/PC-Alex/city_network.net.xml"
        output_file = "sumo_rl/nets/PC-Alex/custom_route_NN.rou.xml"
        fringe_factor = 10
        begin_time = 0
        end_time = 3600

        # Call the function from the other script
        generate_trips_NN.run_full_pipeline(network_file, output_file, fringe_factor, begin_

    def run_simulation(self, traffic_dist, filename):
        """
            Run the simulation and collect (state, reward) tuples
            Param: traffic_dist: intensity of traffic that is used
            Param: filename: name of the save file
        """

        cycle_count = 0
        cycle_length = 90
```

71

```python
 73            # for efficiency have the number of steps be a multiple of the cycle length so that
 74            # every 5 cycle's a data point is taken. Per simulation total_cycles/5 = data points
 75            total_cycles = 30
 76            look_ahead_cycles = 5 # over which horizon the rewards will be calculated
 77            num_steps = total_cycles*cycle_length
 78
 79            # this var keeps track of the rewards over the last 5 cycles
 80            # needs to be reset every simulation, hence not initialized in init
 81            state_buffer = []
 82            reward_buffer = collections.defaultdict(list)
 83
 84            # run the simulation for a amount of steps
 85            for step in range(num_steps):
 86                self.sumo_interface.simulationStep()  # Advances the simulation by one step
 87                self.vehicle_ids = traci.vehicle.getIDList()
 88                # Accumulate the total speed of all vehicles and count the vehicles
 89                speed_all_veh = np.mean([self.sumo_interface.vehicle.getSpeed(vid) for vid in se
 90                arrived_vehicles = self.sumo_interface.simulation.getArrivedIDList()
 91                self.total_arrived_vehicles += len(arrived_vehicles)
 92                self.avg_speeds_per_step.append(np.mean(speed_all_veh))
 93
 94                # keep track of how many cycles are run given a policy
 95                # Every 5 cycles a data point is taken
 96                if step % cycle_length == 0 and step != 0:
 97                    cycle_count += 1
 98                    self.cycle_mean_speed = np.mean(self.avg_speeds_per_step)
 99                    # empty list for next cycle
100                    self.avg_speeds_per_step = []
101
102                    # retrieve the current state
103                    current_state = self.get_state()
104
105                    # retrieve the best tl setting given the state and traffic dist
106                    action = self.get_best_tl_settings(current_state, traffic_dist)
107
108                    # apply the best found tl setting
109                    self.apply_action(action)
110
111                    state_buffer.append(current_state)
112
113                    # Get the reward for this cycle
114                    current_reward = self.calculate_reward()
115
116                    # Add the current reward to all states in the buffer
117                    for state_idx in range(len(state_buffer)):
118                        reward_buffer[state_idx].append(current_reward)
119
120                    first_key = next(iter(reward_buffer))
121
122                    # Check if enough cycles have passed to link a state to its future rewards a
```

72

```python
                if len(reward_buffer[first_key]) == look_ahead_cycles or step + look_ahead_c
                    # Link the oldest state to the cumulative reward of the past x cycles
                    # Initialize sums for each element
                    sum_0, sum_1, sum_2 = 0, 0, 0
                    for x in range(look_ahead_cycles):
                        sum_0 += reward_buffer[first_key][x][0]
                        sum_1 += reward_buffer[first_key][x][1]
                        sum_2 += reward_buffer[first_key][x][2]

                    # Results
                    cumulative_reward = [sum_0, sum_1, sum_2]
                    #print('state buffer pre pop', state_buffer)
                    past_state = state_buffer.pop(first_key) #returns and removes the oldest
                    reward_buffer.pop(first_key)
                    #print('reward_buffer post pop', reward_buffer)
                    self.data.append((past_state, action, cumulative_reward, self.traffic_di
                    self.save_data(filename)  # Save data after each cycle

    def apply_action(self, action):

        """
        Modifies the action order to allign correctly for tl on the perimeter according to t
        Then forwards the action to self.set_red_green_phases to implement the new tl config

        First action value is always inflow/outflow direction of PN forward+right turn
        Second action value is always inflow/outflow direction of PN left turn
        Third action is perpendicular forward+right turn
        Fourth action is perpendicular left turn
        Since the action is given to all tl lights on the perim and the phases always start
        the order must be changed depending on the position of the tl on the perim

        :param action: action given by rl algo
        """
        # Change traffic light configurations based on action
        traffic_lights = self.sumo_interface.trafficlight.getIDList()

        for tl in traffic_lights:
            #checks if the tl is on perimeter, else no modification, SUMO preset phases are
            # line counter is used to determine which line passed the " True" value for the
            on_line, line_counter = self.is_traffic_light_on_perimeter(tl)
            # only modify tl which are on perim line
            if on_line is True:
                #west side tl
                if line_counter == 0:
                    # Create a copy of the action array for the modified action
                    modified_action = action.copy()
                    # Swap the first and third values
                    modified_action[0], modified_action[2] = action[2], action[0]
                    # Swap the second and fourth values
                    modified_action[1], modified_action[3] = action[3], action[1]
```

73

```
173                        # print ('modified action west side tl', modified_action)
174                        self.set_red_green_phases ( modified_action , tl)
175                  # north side tl
176                  if line_counter ==1:
177                        # already correct order
178                        self.set_red_green_phases (action , tl)
179                  # east side tl
180                  if line_counter ==2:
181                        # Create a copy of the action array for the modified action
182                        modified_action = action.copy ()
183                        # Swap the first and third values
184                        modified_action [0] , modified_action [2] = action [2] , action [0]
185                        # Swap the second and fourth values
186                        modified_action [1] , modified_action [3] = action [3] , action [1]
187                        self.set_red_green_phases ( modified_action , tl)
188                  # south side tl
189                  if line_counter ==3:
190                        self.set_red_green_phases (action , tl)
191
192       def get_best_tl_settings (self , state , traffic_dist ):
193             """"
194             This function will use the NN to predict the reward given a TL setting and a state
195             From these the TL setting that gives the best reward will be selected as action
196
197             return: tl setting
198             """"
199             # Load the model from the file
200             # Modify the name of the correct model and scaler that is used based on iteration nu
201             model = load ('iter4_nn_model_rew2.joblib')
202             scaler = load ('x_scaler_rew2_iter4.joblib')
203             best_reward = 0
204             best_tl_setting = []
205
206             def evaluate_expression (traffic_dist ):
207                  """Extract the numerical part from expressions like 'exp(0.0041)', the return wi
208                  return float (traffic_dist [len ("exp("):-1])
209
210             def predict_reward (state , policy , traffic_dist ):
211                  """
212                  Predicts the reward for a given state , policy , and trip count using the trained
213
214                  Parameters:
215                  state (list): The state features as a list of numerical values.
216                  policy (list): The policy features as a list of numerical values.
217                  trip_count (int): The number of trips as an integer.
218
219                  Returns:
220                  float: The predicted reward.
221                  """
222
```

```
223            # modify the format of the traffic dist from string to float
224            traffic_dist = evaluate_expression(traffic_dist)
225
226            # Combine the state, policy, and trip_count into a single feature vector
227            new_input = np.hstack((state, policy, traffic_dist))
228
229            # Standardize the input using the loaded scaler
230            new_input_scaled = scaler.transform([new_input])
231
232            # Predict the reward using the trained model
233            predicted_reward = model.predict(new_input_scaled)
234
235            # Return the predicted reward
236            return predicted_reward
237
238        # List of tl settings for the NN to choose from
239        tl_settings = [
240            [27, 12, 27, 12],
241            [26, 13, 26, 13],
242            [25, 14, 25, 14],
243            [24, 15, 24, 15],
244            [28, 11, 28, 11],
245            [29, 10, 29, 10],
246            [30, 9, 30, 9],
247            [31, 8, 31, 8],
248            [32, 9, 28, 9],
249            [34, 10, 24, 10],
250            [36, 11, 20, 11],
251            [37, 12, 17, 12],
252            [33, 11, 23, 11],
253            [28, 9, 32, 9],
254            [24, 10, 34, 10],
255            [20, 11, 36, 11],
256            [17, 12, 37, 12],
257            [23, 11, 33, 11],
258            [9, 28, 9, 28],
259            [15, 24, 15, 24],
260            [18, 21, 18, 21],
261            [21, 18, 21, 18],
262            [19.5, 19.5, 19.5, 19.5],
263        ]
264
265        list_of_predicted_rewards = []
266        #loop over all tl_settings to find the best one
267        for tl_setting in tl_settings:
268            # Predict the reward
269            predicted_reward = predict_reward(state, tl_setting, traffic_dist)
270            predicted_reward = predicted_reward.item()
271            list_of_predicted_rewards.append(predicted_reward)
272
```

```
273                     # keep track of the best reward and the associated tl_setting
274                     if predicted_reward >= best_reward:
275                         best_tl_setting = tl_setting
276                         best_reward = predicted_reward
277             return best_tl_setting
278
279     def set_simulation_settings(self, policy_settings, traffic_dist_array, route_count_array
280             """
281             Will define the tl policy and traffic and run simulation
282             :param: policy_settings: list of multiple different inputs for the traffic light set
283             :param: traffic_dist_array: list of the type of traffic dist
284             :param: route_count_array: amount of trips that are created
285             """
286             # policy will change the traffic light settings
287             for policy in policy_settings:
288                 # traffic will adjust the kind of traffic in the simulation
289                 for traffic_dist in traffic_dist_array:
290                     # trip count determines how many cars will be in the simulation
291                     for trip_count in route_count_array:
292                         self.traffic_dist = traffic_dist
293                         # keep track of the duration of one simulation
294                         start_time = time.time()
295                         # create the desired traffic for this simulation
296                         self.create_traffic(traffic_dist, trip_count)
297                         self.total_trips = trip_count # keep track of this to use for TCR
298                         # run warm up for each simulation as load-state is used in sumo_cmd
299                         self.sumo_interface.start(self.sumo_cmd)
300                         # adjust the traffic light policy to the used policy after warm up
301                         self.set_traffic_light(policy)
302                         # run the actual simulation
303                         self.run_simulation(traffic_dist, file_name)
304                         self.sumo_interface.close()
305                         print(' *!*!*!*!*!*!*!!*!*!*!*!*!*!*!*!*')
306                         print('------------one simulation duration is', time.time() - start_time
307
308
309     def get_state(self):
310             """
311             Extract the current state
312             Returns: state (array)
313             """
314
315         total_vehicles_inside = 0
316         outside_queue_count = 0
317         inside_queue_count = 0
318
319         # assumed we look at the complete network
320         # elsewise check for controlled tl's -> controlled lanes -> filter for unique
321         travel_times_per_veh = []
322         for vehicle_id in self.vehicle_ids:
```

76

```python
323            travel_times_per_veh.append(self.sumo_interface.vehicle.getAccumulatedWaitingTim
324        total_travel_time = np.sum(travel_times_per_veh)
325        avg_travel_time = np.mean(travel_times_per_veh)
326
327        lane_occupancy_per_lane = [self.sumo_interface.lane.getLastStepOccupancy(lane) for l
328        total_lane_occupancy = sum(lane_occupancy_per_lane)
329        avg_lane_occupancy = np.mean(lane_occupancy_per_lane)
330        vehicle_ids = self.sumo_interface.vehicle.getIDList()
331        total_vehicles = len(vehicle_ids)
332        queu_per_lane = [self.sumo_interface.lane.getLastStepHaltingNumber(lane) for lane in
333        total_queue_length = sum(queu_per_lane)
334        avg_queue_length = np.mean(queu_per_lane)
335        avg_speed = self.cycle_mean_speed
336        inflowing_vehicles, outflowing_vehicles = self.get_inflowing_outflowing_vehicles()
337
338        # Get queue length for each traffic light-controlled lane
339        tl_queue_lengths = []
340        traffic_lights = self.sumo_interface.trafficlight.getIDList()
341        for tl_id in traffic_lights:
342            controlled_lanes = self.sumo_interface.trafficlight.getControlledLanes(tl_id)
343            tl_queue_length = sum(self.sumo_interface.lane.getLastStepHaltingNumber(lane) fo
344            tl_queue_lengths.append(tl_queue_length)
345
346
347        # Calculate the total speed and number of vehicles inside and outside perimeter
348        for vehicle_id in self.vehicle_ids:
349            if self.is_vehicle_inside_perimeter(vehicle_id):
350                total_vehicles_inside +=1
351                if traci.vehicle.getSpeed(vehicle_id) < 0.1:  # Define queue as vehicle movi
352                    inside_queue_count += 1
353            else: #not inside so vehicles outside perimeter
354                if traci.vehicle.getSpeed(vehicle_id) < 0.1:  # Define queue as vehicle movi
355                    outside_queue_count += 1
356        total_vehicles_outside = total_vehicles-total_vehicles_inside
357
358
359        # Combine all state elements including traffic light queue lengths
360        # state = np.array([
361        #     total_lane_occupancy, avg_lane_occupancy, total_vehicles, total_vehicles_inside
362        #     total_vehicles_outside, avg_speed, total_queue_length, avg_queue_length,
363        #     outside_queue_count, inside_queue_count, total_travel_time, avg_travel_time,
364        #     inflowing_vehicles, outflowing_vehicles
365        # ] + tl_queue_lengths)
366
367        state = np.array([total_lane_occupancy, avg_lane_occupancy, total_vehicles, total_ve
    avg_speed,
368                          # 6                                 7                              8
    9                  10                              11
369                          total_queue_length, avg_queue_length, outside_queue_count, inside_q
370                          # 12                              13
```

77

```
371                             inflowing_vehicles , outflowing_vehicles ])
372         return state
373
374     def set_traffic_light ( self , policy ):
375
376             """
377             Modifies the policy order to allign correctly for tl on the perimeter according
378             Then forwards the policy to self.set_red_green_phases to implement the new tl co
379
380             First policy value is always inflow/outflow direction of PN forward+right turn
381             Second policy value is always inflow/outflow direction of PN left turn
382             Third policy is perpendicular forward+right turn
383             Fourth policy is perpendicular left turn
384             Since the policy is given to all tl lights on the perim and the phases always st
385             the order must be changed depending on the position of the tl on the perim
386
387             :param policy: policy given by rl algo
388             """
389             # Change traffic light configurations based on policy
390             traffic_lights = self.sumo_interface.trafficlight.getIDList ()
391
392             for tl in traffic_lights:
393                 #checks if the tl is on perimeter , else no modification , SUMO preset phases
394                 # line counter is used to determine which line passed the " True" value for
395                 on_line , line_counter = self.is_traffic_light_on_perimeter ( tl )
396                 # only modify tl which are on perim line
397                 if on_line is True:
398                     #west side tl
399                     if line_counter == 0:
400                         # Create a copy of the policy array for the modified policy
401                         modified_policy = policy.copy ()
402                         # Swap the first and third values
403                         modified_policy [0] , modified_policy [2] = policy [2] , policy [0]
404                         # Swap the second and fourth values
405                         modified_policy [1] , modified_policy [3] = policy [3] , policy [1]
406                         # print ('modified policy west side tl', modified_policy )
407                         self.set_red_green_phases ( modified_policy , tl )
408                     # north side tl
409                     if line_counter ==1:
410                         # already correct order
411                         self.set_red_green_phases ( policy , tl )
412                     # east side tl
413                     if line_counter ==2:
414                         # Create a copy of the policy array for the modified policy
415                         modified_policy = policy.copy ()
416                         # Swap the first and third values
417                         modified_policy [0] , modified_policy [2] = policy [2] , policy [0]
418                         # Swap the second and fourth values
419                         modified_policy [1] , modified_policy [3] = policy [3] , policy [1]
420                         self.set_red_green_phases ( modified_policy , tl )
```

```
421                     # south side tl
422                     if line_counter ==3:
423                         self.set_red_green_phases(policy, tl)
424
425     def set_red_green_phases(self, action, tl_id):
426         """
427         Applies the action given by the rl algo for all tl on the perimeter (this check is a
428
429         First action value is always inflow/outflow direction of protected network (PN) forw
430         Second action value is always inflow/outflow direction of PN left turn
431         Third action is perpendicular forward+right turn
432         Fourth action is perpendicular left turn
433         Since the action is given to all tl lights on the perim and the phases always start
434         the order must be changed depending on the position of the tl on the perim
435
436         :param: action: action given by rl algo
437         :param: tl_id: traffic light which should be modified
438         """
439
440         # keep track of the original program which will be modified so that we prevent modif
441         if not hasattr(self, 'original_program'):
442             # Get the current (original) program logic
443             current_program_tuple = self.sumo_interface.trafficlight.getAllProgramLogics(tl_
444
445             current_program = current_program_tuple[0]
446
447             # Correct the phase states
448             corrected_phases = self.init_correct_phase_states(current_program.phases)
449
450             # Create a new traffic light program with the corrected phases
451             corrected_program = self.sumo_interface.trafficlight.Logic(
452                 current_program.programID,
453                 current_program.type,
454                 current_program.subParameter,
455                 corrected_phases
456             )
457
458             # Store the modified program
459             # modified is removing the conflicting green phases, action is NOT yet applied
460             self.original_program = corrected_program
461         else:
462             # Use the stored original program for modification
463             current_program = self.original_program
464
465         action_counter = 0
466         modified_phases = []
467         for phase in current_program.phases:
468             # If the phase contains yellow lights ('y'), skip modifying it
469             if 'y' in phase.state or 'Y' in phase.state:
470                 modified_phases.append(phase)  # Keep the yellow phase unchanged
```

79

```
471                     continue
472                 # Modify the duration of the non yellow phases
473                 if action_counter < len(action):
474                     new_duration = action[action_counter]
475                     # prevent durations being smaller than the min green time (should not be pos
476                     if new_duration <= self.min_green_time:
477                         new_duration = self.min_green_time
478                     action_counter += 1
479                 else:
480                     new_duration = phase.duration
481                 # Create a new phase with the modified duration and the same state
482                 modified_phases.append(self.sumo_interface.trafficlight.Phase(new_duration, phas


485         # Create a new traffic light program with the modified phases
486         program = self.sumo_interface.trafficlight.Logic("0", 0, 0, modified_phases)

488         # Set the new traffic light program
489         self.sumo_interface.trafficlight.setProgramLogic(tl_id, program)

491     def calculate_reward(self):
492         """
493         Calculate the reward for the current state-policy pair
494         return: reward (list of 3 floats)
495         """
496         mean_waiting_time_network = 0
497         # Calculate the total speed and number of vehicles inside and outside perimeter
498         vehicle_ids = self.sumo_interface.vehicle.getIDList()
499         for vehicle_id in vehicle_ids:
500             mean_waiting_time_network += self.sumo_interface.vehicle.getWaitingTime(vehicle_

502         # Calculate trip completion rate, total started trips/ arrived trips
503         trip_completion_rate = (self.total_arrived_vehicles / self.total_trips) * 100

505         # play around with the weights for optimal performance and so that all terms have re
506         w1 = 1
507         w2 = 1
508         w3 = 1

510         #reward = -w1*  mean_waiting_time_network + w2 * trip_completion_rate + w3* self.cyc
511         reward = [mean_waiting_time_network, trip_completion_rate,  self.cycle_mean_speed]


514         return reward

516     def save_data(self, filename):
517         """Append collected data to a CSV file continuously."""
518         df = pd.DataFrame(self.data, columns=["state", "policy", "reward", "traffic_dist", "
519         # Use 'mode="a"' to append, and 'header=False' to avoid writing the header repeatedl
520         df.to_csv(filename, mode='a', header=False, index=False)
```

80

```
521             # Clear data after saving to prevent duplicate entries
522             self.data = []
523
524     def is_traffic_light_on_perimeter(self, tl_id):
525             """
526             Checks if a traffic light is on the perimeter
527             :param tl_id: The ID of the traffic light.
528             :return: True if the tl is on the perimeter, False otherwise.
529
530             Note that this approach does not work for non-straight roads
531
532             """
533             lane_coordinates = []
534             x_coordinates = []
535             y_coordinates = []
536
537             # Get a list of lanes controlled by the tl
538             controlled_lanes = traci.trafficlight.getControlledLanes(tl_id)
539
540             # Loop over all lanes and get their start and end coordinates. (all lanes are straig
541             for lane in controlled_lanes:
542                 lane_coordinates.append(traci.lane.getShape(lane))
543
544             # Get the x and y coordinates of the end points of the lane (where the tl is),
545             for x in range(len(controlled_lanes)):
546                 x_coordinates.append(lane_coordinates[x][-1][0])
547                 y_coordinates.append(lane_coordinates[x][-1][1])
548
549             x_coordinate_tl = np.mean(x_coordinates)
550             y_coordinate_tl = np.mean(y_coordinates)
551
552
553             xmin, ymin, xmax, ymax = self.perim_box
554             coordinates_tl = [x_coordinate_tl, y_coordinate_tl]
555
556             # Define the corners of the perimeter
557             corners = [
558                 (xmin, ymin),   # Bottom-left corner
559                 (xmin, ymax),   # Top-left corner
560                 (xmax, ymax),   # Top-right corner
561                 (xmax, ymin)    # Bottom-right corner
562             ]
563
564             # Check if the TL is located at any of the corners, corners will not be considered a
565             for corner in corners:
566                 if np.isclose(coordinates_tl[0], corner[0], atol=10) and np.isclose(coordinates_
567                     # print('tl id at corner is', tl_id)
568                     return False, None  # TL is at the corner, return False
569
570             # Define the sides of the perimeter box as LineString objects
```

```
571          lines = [
572              LineString([(xmin, ymin), (xmin, ymax)]),   # Left side
573              LineString([(xmin, ymax), (xmax, ymax)]),   # Top side
574              LineString([(xmax, ymax), (xmax, ymin)]),   # Right side
575              LineString([(xmax, ymin), (xmin, ymin)])    # Bottom side
576          ]
577
578
579          # Depending on which iteration of line the statement is true, the position of the tl
580          # Each loop of the line in lines correspond to the next lines above so left, top, ri
581          # Check if the point is near any of the lines
582          line_counter = 0
583          for line in lines:
584              if self.point_near_line(coordinates_tl, line):
585                  on_line = True
586                  return on_line, line_counter
587              else:
588                  on_line = False
589                  line_counter += 1
590
591          return on_line, line_counter
592
593      def is_vehicle_inside_perimeter(self, vehicle_id):
594          """
595          Check if a vehicle is inside the perimeter.
596
597          :param vehicle_id: The ID of the vehicle.
598          :return: True if the vehicle is inside the perimeter, False otherwise.
599          """
600          # x, y coordinate of vehicle
601          x, y = traci.vehicle.getPosition(vehicle_id)
602          # coordinates of the perimeter
603          xmin, ymin, xmax, ymax = self.perim_box
604          # check if a vehicle is within the perimeter
605          return xmin <= x <= xmax and ymin <= y <= ymax
606
607      def point_near_line(self, point, line, tolerance=20.0):
608          """
609          Checks if the point is within a certain tolerance of a line.
610
611          :param point: A tuple (x, y) representing the point.
612          :param line: A LineString object representing the line.
613          :param tolerance: The tolerance within which the point should be from the line (in m
614          :return: True if the point is within the tolerance of the line, False otherwise.
615          """
616          point_obj = Point(point)
617          # determines the euclidian distance from the point to the line
618          distance = point_obj.distance(line)
619          # If the distance is more than the tolerance it will return False, else True
620          return distance <= tolerance
```

82

```
621
622     def get_inflowing_outflowing_vehicles(self):
623         """
624             Get the count of vehicles inflowing into the perimeter and outflowing out of the
625             Return: number of inflowing vehicles (int)
626             Return: number of outflowing vehicles (int)
627         """
628         inflow_count = 0
629         outflow_count = 0
630
631         # Track the vehicles' previous positions (inside or outside perimeter)
632         if not hasattr(self, 'previous_positions'):
633             self.previous_positions = {}
634
635         for vehicle_id in self.vehicle_ids:
636             # Check if the vehicle is currently inside the perimeter
637             is_inside_now = self.is_vehicle_inside_perimeter(vehicle_id)
638             was_inside_before = self.previous_positions.get(vehicle_id, False)
639
640             # Count inflow if the vehicle was outside before but is now inside
641             if not was_inside_before and is_inside_now:
642                 inflow_count += 1
643             # Count outflow if the vehicle was inside before but is now outside
644             if was_inside_before and not is_inside_now:
645                 outflow_count += 1
646
647             # Update the previous position tracking
648             self.previous_positions[vehicle_id] = is_inside_now
649
650         return inflow_count, outflow_count
651
652
653     def init_correct_phase_states(self, phases):
654         """
655         Modifies the states of all traffic light phases by changing all 'g' (green with
656         We only want G (green without confliction)
657         Parameters:
658         - phases (list of self.sumo_interface.trafficlight.Phase): A list of traffic lig
659
660         Returns:
661         - list of self.sumo_interface.trafficlight.Phase: The modified list of traffic l
662         """
663
664         # all conflicted greens should be removed already in netedit
665         corrected_phases = []
666         for phase in phases:
667             # Modify the phase state to change all 'g' to 'r'
668             new_state = phase.state.replace('g', 'r')
669             # Create a new phase with the modified state and original duration
670             corrected_phase = self.sumo_interface.trafficlight.Phase(phase.duration, new
```

```
671                      corrected_phases . append ( corrected_phase )
672              return corrected_phases
673
674   sumoBinary = sumolib . checkBinary ( 'sumo ') # change to sumo - gui if you want to use gui . gui ca
675   sumo_cmd = [
676       sumoBinary ,
677       "-n", "sumo_rl/nets/PC-Alex/City_network.net.xml",
678       "-r", "sumo_rl/nets/PC-Alex/custom_route_NN.rou.xml" ,
679       "--fcd-output", "fcd.xml",
680       "--vehroute-output", "vehroutes.xml",
681       "--default.speeddev", "0",
682       "--no-step-log",
683       "--summary-output", "summary_output.xml",
684       "--load-state", "City_warm_up.xml" # using load_state with city warm up directly uses th
685       # state of city_warm_up as starting point of the simulation which means no additional wa
686    ]
687
688
689   data_collector = SumoDataCollector ( sumo_cmd )
690
691
692   def initialize_settings ( run_training = True ):
693       """
694       Initialize_settings for either training or validation based on the input .
695
696       Parameters :
697       run_training ( bool ): If True , run training simulations . If False , run validation simulat
698       """
699       if run_training :
700           # Training settings
701           #phase  = 90s , yellow time = 3*4=12 , remaining time 78
702
703           policy_settings = [
704               [29, 10, 29, 10],  # standard traffic light configurations
705               [31, 8, 31, 8],    # extended duration for main phases
706               [27, 12, 27, 12],  # extended duration for the secondary phases
707               [32, 11, 24, 11],  # extended duration for phase into perimeter
708               [24, 11, 32, 11],  # extended duration for phase out of perimeter
709               [30, 11, 26, 11],  # equally divided phase durations
710               [28, 11, 28, 11],  # standard traffic light configurations
711               [31, 8, 31, 8],    # extended duration for main phases
712               [25, 14, 25, 14],  # extended duration for the secondary phases
713               [35, 11, 21, 11],  # extended duration for phase into perimeter
714               [21, 11, 35, 11],  # extended duration for phase out of perimeter
715               [19.5, 19.5, 19.5, 19.5],
716               [26, 13, 26, 13],
717               [29, 10, 26, 13],
718               [28, 13, 25, 12],
719               [34, 9, 25, 10],
720               [33, 9, 26, 10],
```

84

```
721              [32, 10, 26, 10],
722              [30, 12, 26, 10],
723              [27, 11, 27, 13],
724              [25, 11, 30, 12],
725              [28, 12, 26, 12],
726              [33, 9, 24, 12],
727              [32, 9, 25, 12],
728              [31, 10, 25, 12],
729              [26.5, 12.5, 26.5, 12.5],
730
731          ]
732          traffic_dist_array = ['exp(0.0041)','exp(0.00435)', 'exp(0.00445)', 'exp(0.00425)',
733                                'exp(0.0044)','exp(0.0043)', 'exp(0.0045)', 'exp(0.0046)', 'e
734                                'exp(0.00475)', 'exp(0.00485)', 'exp(0.00415)',
     'exp(0.00435)']
735          route_count_array = [1000] # this is kept constant, instead traffic intensity is cha
736          # if this is changed as well the traffic intensity changes too drastically
737      else:
738          # Validation settings
739          policy_settings = [
740              [20, 15, 20, 23],
741              [18, 18, 21, 21],
742              [17, 19, 20, 22],
743              [24, 12, 18, 24],
744              [19, 20, 19, 20],
745              [15, 17, 23, 23],
746              [21, 11, 23, 23],
747              [22, 14, 19, 23],
748              [16, 18, 22, 22],
749              [23, 9, 23, 23]
750          ]
751          traffic_dist_array = ['exp(0.0047)', 'exp(0.0042)','exp(0.0049)', 'exp(0.00428', 'ex
752          route_count_array = [1000]
753      return policy_settings, traffic_dist_array, route_count_array
754
755
756  run_training = True # True means training data will be created, false means validation data
757
758  # modify the name of the save file to match the iter run
759  if run_training is True:
760      file_name = 'sumo_training_data_iter5.csv'
761  else:
762      file_name = 'sumo_validation_data_iter5.csv'
763
764  policy_settings, traffic_dist_array, route_count_array = initialize_settings(run_training)
765
766  data_collector.set_simulation_settings(policy_settings, traffic_dist_array, route_count_array
```

```
1  import matplotlib.pyplot as plt
```

```
2   from sklearn.neural_network import MLPRegressor
3   from sklearn.model_selection import train_test_split, GridSearchCV, RandomizedSearchCV, cros
4   from sklearn.preprocessing import MinMaxScaler
5   from sklearn.metrics import mean_absolute_error
6   import numpy as np
7   import pandas as pd
8   import ast
9   import re
10  from joblib import dump
11
12  # Load and preprocess the data
13  train_filename = "sumo_training_data2.csv"
14  train_df = pd.read_csv(train_filename, names=["state", "policy", "reward", "traffic_dist", "
15
16  valid_filename = "sumo_validation_data_testing.csv"
17  valid_df = pd.read_csv(valid_filename, names=["state", "policy", "reward", "traffic_dist", "
18
19
20  # Function to format and convert state vector strings
21  def format_state_vector(state_str):
22      # Replace spaces with commas while ignoring scientific notation parts
23      state = re.sub(r'(?<=\d)\s+(?=[\d\-])', ', ', state_str)
24      return state
25  train_df["state"] = train_df["state"].apply(format_state_vector)
26  valid_df["state"] = train_df["state"].apply(format_state_vector)
27
28  # Define various subsets of state vector indices
29  state_index_combinations = [
30      [0,2,3,4,5,6,8,9,10,12,13], #avg
31      [1,2,3,4,5,7,8,9,11,12,13], #total
32      [1, 3, 4, 5, 8, 9, 12, 13], # intuitively most important
33      [5], # only avg speed
34      [5,8,9], # avg speed, queue inside and outside
35      [3,4,5], # vehicles inside, outside, avg speed
36  ]
37
38  def preprocess_data(df, state_indices=None):
39      """ Preprocess data with optional state vector subset selection. """
40      def fix_format(data_string):
41          data_string = data_string.strip()
42          data_string = re.sub(r'\s+', ' ', data_string)
43          fixed_string = re.sub(r'(\d)\s+(\d)', r'\1, \2', data_string)
44          return fixed_string
45
46      df["state"] = df["state"].apply(fix_format).apply(ast.literal_eval)
47      df["policy"] = df["policy"].apply(fix_format).apply(ast.literal_eval)
48
49      def extract_inner_value(value):
50          return float(value[len("exp("):-1])
51      df['traffic_dist'] = df['traffic_dist'].apply(extract_inner_value)
```

```
52
53
54    states = df["state"].apply(lambda x: [x[i] for i in state_indices] if state_indices else
55    policies = df["policy"].tolist()
56    rewards = df["reward"].tolist()
57    traffic_dist = np.array(df["traffic_dist"].tolist()).reshape(-1, 1)
58
59    X = np.hstack((states, policies, traffic_dist))
60    reward_number = 2
61    y = np.array([ast.literal_eval(reward)[reward_number] for reward in rewards], dtype=np.f
62
63    return X, y
64
65  # Set up scalers
66  x_scaler = MinMaxScaler()
67  y_scaler = MinMaxScaler()
68
69  def run_optimize_hyperparameter(X_train, y_train, X_valid, y_valid):
70      """ Optimize the NN hyperparameters with random search and return best model and perform
71      model = MLPRegressor(max_iter=5000, random_state=42, early_stopping=True)
72
73      # Define the parameter distribution for Random Search
74      param_dist = {
75          'hidden_layer_sizes': [(128, 128, 64, 32), (128, 64, 32, 16), (256, 128, 64, 32, 16)
76          'activation': ['relu'],
77          'solver': ['adam'],
78          'alpha': np.logspace(-4, -1, 4),
79          'learning_rate': ['constant', 'adaptive'],
80          'batch_size': [32, 64, 128],
81          'early_stopping': [True]
82      }
83
84      # Set up Random Search
85      random_search = RandomizedSearchCV(
86          estimator=model,
87          param_distributions=param_dist,
88          scoring='neg_mean_squared_error',
89          cv=3,
90          n_iter=200,  # Number of random combinations to try
91          verbose=2,
92          random_state=42,
93          n_jobs=-1
94      )
95
96      random_search.fit(X_train, y_train)
97
98      # Get the best model and evaluate on validation set
99      best_model = random_search.best_estimator_
100     best_model.fit(X_train, y_train)  # Fit best model on full training data
101
```

```python
102        predictions = best_model.predict(X_valid)
103        original_scale_predictions = y_scaler.inverse_transform(predictions.reshape(-1, 1)).flat
104        original_scale_y_valid = y_scaler.inverse_transform(y_valid.reshape(-1, 1)).flatten()
105
106        mae = mean_absolute_error(original_scale_y_valid, original_scale_predictions)
107        rmse = np.sqrt(np.mean((original_scale_y_valid - original_scale_predictions) ** 2))
108
109        return random_search.best_params_, mae, rmse
110
111 # Experiment with different state subsets
112 results = []
113 for state_indices in state_index_combinations:
114     print(f"Testing state indices: {state_indices}")
115
116     # Preprocess data with the selected subset of state indices
117     X_train, y_train = preprocess_data(train_df, state_indices)
118     X_valid, y_valid = preprocess_data(valid_df, state_indices)
119
120     # Scale the data
121     X_train = x_scaler.fit_transform(X_train)
122     X_valid = x_scaler.transform(X_valid)
123     y_train = y_scaler.fit_transform(y_train.reshape(-1, 1)).flatten()
124     y_valid = y_scaler.transform(y_valid.reshape(-1, 1)).flatten()
125
126     # Run NN training and get performance metrics
127     best_params, mae, rmse = run_optimize_hyperparameter(X_train, y_train, X_valid, y_valid)
128     results.append({
129         'state_indices': state_indices,
130         'MAE': mae,
131         'RMSE': rmse,
132         'best_params': best_params
133     })
134     # print(f"MAE: {mae}, RMSE: {rmse}")
135
136
137 # Convert results to a DataFrame and display
138 results_df = pd.DataFrame(results)
139
140 # Find the top 3 rows with the lowest MAE
141 top_3_mae = results_df.nsmallest(3, 'MAE')[['state_indices', 'MAE', 'RMSE', 'best_params']]
142
143 # Find the top 3 rows with the lowest RMSE
144 top_3_rmse = results_df.nsmallest(3, 'RMSE')[['state_indices', 'MAE', 'RMSE', 'best_params']
145
146 # Save the top 3 MAE and top 3 RMSE results to CSV files
147 top_3_mae.to_csv("top_3_mae_results.csv", index=False)
148 top_3_rmse.to_csv("top_3_rmse_results.csv", index=False)
149
150 # Save the results DataFrame to a CSV file
151 results_df.to_csv("nn_optimization_results.csv", index=False)
```

```
152
153  # Display the results
154  print("Top 3 by MAE:")
155  print(top_3_mae)
156
157  print("\nTop 3 by RMSE:")
158  print(top_3_rmse)
159
160  print(results_df)
```

```
1   import gymnasium
2   import gymnasium.spaces.discrete
3   import traci
4   import numpy as np
5   from stable_baselines3 import PPO
6   import os
7   import time
8   os.environ['TF_CPP_MIN_LOG_LEVEL'] = '2'  # Suppress INFO and WARNING logs
9   import tensorflow as tf
10  from stable_baselines3.common.callbacks import BaseCallback
11  tf.debugging.set_log_device_placement(True)
12  import matplotlib.pyplot as plt
13  from shapely.geometry import Point, LineString
14  from stable_baselines3.common.vec_env import SubprocVecEnv
15  from stable_baselines3.common.monitor import Monitor
16  import sumolib
17  import re
18  import generate_trips_NN
19  # import uuid
20  #import self.sumo_interface as tc
21
22
23  # class AvgSpeedLoggingCallback(BaseCallback):
24  #     def __init__(self, summary_writer, verbose=0):
25  #         super(AvgSpeedLoggingCallback, self).__init__(verbose)
26  #         self.summary_writer = summary_writer
27
28  #     def _on_step(self):
29
30  #         avg_speeds = []
31  #         for env in self.model.get_env().envs:  # Get all environments
32  #             avg_speeds.append(env.sumo_env.self.avg_speed)  # Assuming `avg_speed` attribu
33
34  #         # Calculate the average speed across all environments
35  #         avg_speed_across_envs = sum(avg_speeds) / len(avg_speeds)
36
37  #         # Log the average speed into TensorBoard using the provided summary writer
38  #         with self.summary_writer.as_default():
39  #             tf.summary.scalar('avg_speed_across_envs', avg_speed_across_envs, step=self.ep
```

89

```python
40
41
42  class SumoEnvironment ():
43      """
44          Environment for SUMO
45          Class that contains all functions that interact with SUMO
46      """
47
48      def __init__(self, sumo_cfg, max_steps, rank, use_libsumo=False):
49
50          # Choose between sumo and libsumo
51          self.use_libsumo = use_libsumo
52
53          # Import the correct module based on `use_libsumo`
54          if self.use_libsumo:
55              import libsumo as sumo_interface
56              print(' using libsumo')
57          else:
58              import traci as sumo_interface
59              print(' using traci')
60          self.sumo_interface = sumo_interface
61
62          # main config file provided
63          self.sumo_cfg = sumo_cfg
64          # max seconds of simulation
65          self.max_steps = max_steps
66          self.current_step = 0
67          self.avg_speeds_per_step = []
68          self.all_avg_speeds = []
69          self.use_gui = False
70          self.episode_counter = 0
71          self.phases = []
72          self.avg_speed_per_episode = []
73          self.previous_phase = {}
74          self.log_file_path = "avg_speed_log_phases_simple.txt"  # Path to the log file
75          # xmin, ymin, xmax, ymax coordinates of the perimeter in meters
76          self.perim_box = 0, -200, 800, 600 #modify this if the permiter changes
77          self.num_rewards = 0
78          self.reward_mean = 0
79          self.reward_var = 0
80          self.avg_speed_network = 0
81          self.min_green_time = 8
82          self.on_line_tl_ids = []
83          self.avg_speed_over_simulation = []
84          self.episode_90_sec_mean = []
85          self.rank = rank
86          self.vehicle_data = []
87          self.avg_speed = []
88          self.avg_speed_90sec = []
89          self.avg_speed_per_sec = []
```

90

```
90          self.mean_waiting_time_90_sec = []
91          self.mean_waiting_per_sec = []
92          self.waiting_time_per_second = []
93          self.change_in_waiting_time_per_second = []
94          self.reward_list = []
95          self.model = None #later defined in set_model
96          self.prev_waiting_time = 0
97          self.change_in_waiting_time = []
98          self.speed_inside_perim_90_sec = []
99          self.speed_outside_perim_90_sec = []
100         self.avg_speed_inside_perim_per_sec = []
101         self.avg_speed_outside_perim_per_sec = []
102         self.avg_speed_inside_perim = []
103         self.avg_speed_outside_perim = []
104         self.start_time = 0
105         self.end_time = 0
106         self.total_arrived_vehicles = 0
107         self.total_departed_vehicles = 0
108         self.traffic_light_data = {}
109         self.traffic_lights = []
110         self.times_reset_is_called = 0
111         #self.create_traffic()
112
113         super(SumoEnvironment, self).__init__()
114
115         # Initialize log file
116         with open(self.log_file_path, 'w') as f:
117             f.write("Episode,Average Speed\n")  # Write header
118
119
120     def set_model(self, model):
121         """Sets the model to the environment, so that we can log PPO-related information in
122         self.model = model
123
124     @staticmethod
125     def make_sumo_env(sumo_cfg, max_steps, rank, seed=0):
126         """
127         Utility function to create a sumo environment and set the seed for it.
128         :param sumo_cfg: SUMO config file path
129         :param max_steps: maximum steps per episode
130         :param rank: rank of the environment (index)
131         :param seed: random seed for the environment
132         :return: the environment instance
133         """
134         def _init():
135             env = SumoGymEnv(sumo_cfg, max_steps, rank)
136             env = Monitor(env)  # Wrapping the environment with Monitor
137             env.rank = rank
138             env.reset(seed=seed + rank)  # Set a different seed for each environment
139             return env
```

91

```
140             return _init
141
142
143     def reset(self, *, seed=None, options=None):
144         """
145         Reset function that resets the SUMO environment
146         ****implement starting sim x amount of steps in, add seed, options******
147
148         Parameters:
149         - seed: optional specific seed for the simulation run
150         - options:
151
152         Returns:
153         - state: state of the environment
154         - info: mandatory param for RL algo, can be empty dict
155         """
156         print(f"----------------Reset called for episode {self.episode_counter}")
157         # if the reset function is called before the start of the sim self.sumo_interface co
158         # if reset is called at the end of an episode the self.sumo_interface connection mus
159         if self.episode_counter >= 1:
160             # this loop only enters after first episode is completed, and thus start_time is
161             # var used before defined can thus be ignored
162             #print('episode duration is', time.time() - self.start_time)
163             self.sumo_interface.close()
164         # else:
165         #     """For more stochasticity in the simulation create traffic could be called eac
166         #     self.create_traffic()
167
168         if self.current_step == 0:
169             self.start_time = time.time()
170
171         #reset variables:
172         self.total_arrived_vehicles = 0
173         self.total_departed_vehicles = 0
174         self.current_step = 0
175
176         #self.start_time = time.time
177         # Run the simulation with or without gui depending on the value of self.use_gui whic
178         # sumo_binary = "sumo-gui" if self.use_gui else "sumo"
179         # if sumo_binary == "sumo-gui" and self.use_libsumo == True:
180         #     print(' no gui possible with libsumo')
181
182         # Build the command to start SUMO
183         # sumo_cmd = [sumo_binary, "-c", self.sumo_cfg, "--no-warnings", "true", "--log", "l
184
185
186         print(' new routing used')
187         self.times_reset_is_called +=1
188         # print(' episode couner', self.episode_counter)
189         # uncomment "-r", "sumo_rl/nets/PC-Alex/unique_routes.rou.xml", to use original rout
```

```
190
191         # used for libumo but does not seem to be faster
192         sumoBinary = sumolib.checkBinary('sumo-gui') if self.use_gui else  sumolib.checkBina
193         sumo_cmd = [
194             sumoBinary,
195             "-n", "sumo_rl/nets/PC-Alex/City_network.net.xml",
196             "-r", "sumo_rl/nets/PC-Alex/custom_route_NN.rou.xml",
197             "--fcd-output", "fcd.xml",
198             "--vehroute-output", "vehroutes.xml",
199             "--default.speeddev", "0",
200             "--no-step-log",
201             "--load-state", "City_warm_up.xml"]
202
203         # Set random seed if provided
204         if seed is not None:
205             sumo_cmd.extend(["--seed", str(seed)])  # Add seed option to the SUMO command
206
207         # Start SUMO
208         self.sumo_interface.start(sumo_cmd)
209
210         if self.episode_counter == 0:
211             self.traffic_lights = self.sumo_interface.trafficlight.getIDList()
212             for tl_id in self.traffic_lights:
213                 on_line, line_counter = self.is_traffic_light_on_perimeter(tl_id)
214                 if on_line == True:
215                     # Store the data of tl on perim in the dictionary with tl_id as the key
216                     self.traffic_light_data[tl_id] = {
217                         'line_counter': line_counter
218                     }
219
220     """"Warm up np longer required, it is now loaded in using sumo_cmd load"""
221
222         # Initialize the current step and increase the episode counter
223         self.current_step = 0
224         self.episode_counter += 1
225
226         # Get the state of the environment
227         state = self.get_state()
228
229         # Info is a mandatory return argument for PPO and STB3, an empty dict is sufficicent
230         info = {}
231         return state, info
232
233     def create_traffic(self):
234         """"
235         Generates the traffic flows for the simulation.
236         param: trip_count: how many trips are present during the simulation (how busy the ne
237         param: fringe_factor: The ratio between traffic spawning at the outside of the simul
238
239         """"
```

```
240
241         # Now call the function as needed in the current script
242         network_file = "sumo_rl/nets/PC-Alex/city_network.net.xml"
243         output_file = "sumo_rl/nets/PC-Alex/custom_route_NN.rou.xml"
244         fringe_factor = 1
245         begin_time = 0
246         end_time = 3600
247         traffic_dist = 'exp(0.0044)'
248         trip_count = 1000
249
250         print(' ----------------------output file is', output_file)
251         # Call the function from the other script
252         generate_trips_NN.run_full_pipeline(network_file, output_file, fringe_factor, begin_
253
254     def step(self, action):
255         """
256         Reset function that resets the SUMO environment
257         ***Remove the additional way to get avg speed after error testing*****
258
259         Parameters:
260         - seed: optional specific seed for the simulation run
261         - options:
262
263         Returns:
264         - state: state of the environment
265         - info: mandatory param for RL algo, can be empty dict
266         """
267
268         # adjust the action such that the cycle time remains constant
269         # Base phase durations (without actions applied)
270         # Hard coded, should change if netedit durations are changed or extract from sumo
271
272         if np.isnan(action).any():
273             print(' nan in action space')
274         phase_durations = np.array([28, 11, 28, 11])
275
276         # # Compute the adjusted phase durations using the action values
277         adjusted_durations = phase_durations * action
278
279         # # Compute the sum of the adjusted phase durations
280         sum_adjusted_durations = np.sum(adjusted_durations)
281
282         # # total sum of phase durations (cycle time is sum_of_phase_durations + 4*3 (yellow
283         sum_of_phase_durations = sum(phase_durations)
284
285         # # Scale the actions so that the sum of adjusted durations equals 78 seconds
286         scaled_actions = action * (sum_of_phase_durations / sum_adjusted_durations)
287
288         # Apply action
289         self.apply_action(scaled_actions)
```

```
290
291          # Run simulation for 90 steps
292          timestep = 90
293          start_time_90_sec_loop = time.time()
294          done = False
295          for _ in range(timestep):
296              self.sumo_interface.simulationStep()
297              # stop the simulation if it is empty
298              if self.sumo_interface.vehicle.getIDCount() == 0:
299                  done = True
300                  print('------steps required to stop is', self.current_step)
301                  break
302
303              self.current_step += 1
304              arrived_vehicles = self.sumo_interface.simulation.getArrivedIDList()
305              self.total_arrived_vehicles += len(arrived_vehicles)
306              departed_vehicles = self.sumo_interface.simulation.getDepartedIDList()
307              self.total_departed_vehicles += len(departed_vehicles)
308
309              # # Get newly departed and arrived vehicles
310              # new_vehicles = self.sumo_interface.simulation.getDepartedIDList()
311              # print(' new vehicles are', new_vehicles)
312              # arrived_vehicles = self.sumo_interface.simulation.getArrivedIDList()
313              # print('arrived vehicles are', arrived_vehicles)
314
315              # # Subscribes speed and position for new vehicles
316              # for vehicle_id in new_vehicles:
317              #     print(' new vehicle id is', vehicle_id)
318              #     print(' tc.var speed is', tc.VAR_SPEED)
319              #     print('tc.VAR_POSITION' , tc.VAR_POSITION)
320              #     self.sumo_interface.vehicle.subscribe(vehicle_id, [tc.VAR_SPEED, tc.VAR_PO
321              #     print(' subscribe check', self.sumo_interface.vehicle.getSubscriptionResul
322
323              # # Unsubscribe from vehicles that have arrived
324              # for vehicle_id in arrived_vehicles:
325              #     # Check if the vehicle has an active subscription before attempting to uns
326              #     subscription_results = self.sumo_interface.vehicle.getSubscriptionResults(
327              #     if subscription_results:
328              #         print('if statement test')
329              #         self.sumo_interface.vehicle.unsubscribe(vehicle_id)
330
331              # # Retrieve and use the subscribed data for all active vSehicles
332              # for vehicle_id in new_vehicles:
333              #     vehicle_data = self.sumo_interface.vehicle.getSubscriptionResults(vehicle_
334              #     if vehicle_data:
335              #         speed = vehicle_data[tc.VAR_SPEED]
336              #         # keep track of the speed of all cars
337              #         avg_speed_array.append(speed)
338              # # save the avg speed over all cars for each time step
339              # self.avg_speed_90sec.append(np.mean(avg_speed_array))
```

95

```
340
341            vehicle_ids = self.sumo_interface.vehicle.getIDList()
342            speed_per_vehicle = []
343            waiting_time_vehicle = []
344
345            car_id_inside_perim = []
346            car_id_outside_perim = []
347            speed_inside_perim = []
348            speed_outside_perim = []
349            # Accumulate the total speed of all vehicles and count the vehicles
350            for veh_id in vehicle_ids:
351                # speed of all vehicles
352                speed_per_vehicle.append(self.sumo_interface.vehicle.getSpeed(veh_id))
353                waiting_time_vehicle.append(self.sumo_interface.vehicle.getWaitingTime(veh_i
354                if self.is_vehicle_inside_perimeter(veh_id):
355                    car_id_inside_perim.append(veh_id)
356                    speed_inside_perim.append(self.sumo_interface.vehicle.getSpeed(veh_id))
357                else:
358                    car_id_outside_perim.append(veh_id)
359                    speed_outside_perim.append(self.sumo_interface.vehicle.getSpeed(veh_id))
360
361            # Save the avg waiting time per second
362            self.waiting_time_per_second.append(np.mean(waiting_time_vehicle))
363
364            # mean of all vehicles per sec
365            self.avg_speed_per_sec.append(np.mean(speed_per_vehicle))
366
367            self.mean_waiting_per_sec.append(np.mean(waiting_time_vehicle))
368
369            if len(speed_inside_perim) >0:
370                self.avg_speed_inside_perim_per_sec.append(np.mean(speed_inside_perim))
371            else:
372                self.avg_speed_inside_perim_per_sec.append(0)
373
374            if len(speed_outside_perim) >0:
375                self.avg_speed_outside_perim_per_sec.append(np.mean(speed_outside_perim))
376            else:
377                self.avg_speed_outside_perim_per_sec.append(0)
378
379
380        # ***outside the 90 sec loop ****
381        for tl_id, data in self.traffic_light_data.items():
382                line_counter = data['line_counter']  # Extract the line_counter array
383                controlled_lanes = self.sumo_interface.trafficlight.getControlledLanes(tl_id
384                # for some reason SUMO takes duplicates, these are removed by conversion to
385                unique_controlled_lanes = tuple(set(controlled_lanes))
386                # the lanes should also be order in a logical order so that the correct lane
387                sorted_lanes = self.sort_lanes_clockwise(unique_controlled_lanes)
388                #west side tl
389                if line_counter == 0:
```

```
390                    vehicles_on_entering_lane = []
391                    vehicles_on_sec_link = []
392
393                    # phases are as follows , starting at north and going clockwise:
394                    # north 1,2,3
395                    # east 3,4,5
396                    # south 6,7,8
397                    # west 9, 10, 11
398
399                    # controlled lanes are slightly different:
400                    # north 1,2
401                    # east 3,4
402                    # south 5,6
403                    # west 7,8
404
405                    # save the lanes which go into the perim in the tl data dict if it's not
406                    # does not change over the duration of the sim, so must only be filled o
407
408                    # fill up the dictionary 'self.traffic_light_data[]
409                    # dict will contain: {line_counter , lanes_into_perim , vehicles_on_enteri
410                    if 'lanes_into_perim' not in self.traffic_light_data[tl_id] or not self.
411                        self.traffic_light_data[tl_id].setdefault('lanes_into_perim', [])
412                    lanes_into_perim = [sorted_lanes[4], sorted_lanes[5]]
413                    self.traffic_light_data[tl_id]['lanes_into_perim'].extend(lanes_into_per
414
415                    # Loop through each lane in lanes_into_perim and get the vehicles that e
416                    for lane in lanes_into_perim:
417                        if len(self.sumo_interface.lane.getLastStepVehicleIDs(lane)) > 0:
418                            vehicles_on_entering_lane.extend(self.sumo_interface.lane.getLas
419
420                    # save the vehicle ids of the vehicles on entering and sec lanes
421                    # key only needs to be created once
422                    if 'vehicles_on_entering_lane' not in self.traffic_light_data[tl_id] or
423                        self.traffic_light_data[tl_id].setdefault('vehicles_on_entering_lane
424                    self.traffic_light_data[tl_id]['vehicles_on_entering_lane'].extend(vehic
425
426                    # loop over all veh on the entering lane
427                    inflow_enter = 0
428                    queue_enter = 0
429                    for vehid in vehicles_on_entering_lane:
430                        # check which ones are not queueing
431                        if self.sumo_interface.vehicle.getSpeed(vehid) > 0.1:
432                            inflow_enter += 1
433                        # remainder of cars are queuing
434                        else:
435                            queue_enter += 1
436                    # save the amount of vehicles that queue to enter the pc and on the seco
437                    # key only needs to be created once
438                    if 'queue_enter' not in self.traffic_light_data[tl_id] or not self.traff
439                        self.traffic_light_data[tl_id].setdefault('queue_enter', [])
```

97

```
440                    self.traffic_light_data[tl_id].setdefault('inflow_enter', [])
441                self.traffic_light_data[tl_id]['queue_enter'] = queue_enter
442                self.traffic_light_data[tl_id]['inflow_enter'] = inflow_enter
443
444            #north side tl
445            if line_counter == 1:
446                vehicles_on_entering_lane = []
447                vehicles_on_sec_link = []
448                # fill up the dictionary 'self.traffic_light_data[]
449                # dict will contain: {line_counter, lanes_into_perim, vehicles_on_enteri
450                if 'lanes_into_perim' not in self.traffic_light_data[tl_id] or not self.
451                    self.traffic_light_data[tl_id].setdefault('lanes_into_perim', [])
452                lanes_into_perim = [sorted_lanes[6], sorted_lanes[7]]
453                self.traffic_light_data[tl_id]['lanes_into_perim'].extend(lanes_into_per
454
455                # Loop through each lane in lanes_into_perim and get the vehicles that e
456                for lane in lanes_into_perim:
457                    if len(self.sumo_interface.lane.getLastStepVehicleIDs(lane)) > 0:
458                        vehicles_on_entering_lane.extend(self.sumo_interface.lane.getLas
459
460                # save the vehicle ids of the vehicles on entering and sec lanes
461                # key only needs to be created once
462                if 'vehicles_on_entering_lane' not in self.traffic_light_data[tl_id] or
463                    self.traffic_light_data[tl_id].setdefault('vehicles_on_entering_lane
464                self.traffic_light_data[tl_id]['vehicles_on_entering_lane'].extend(vehic
465
466                # loop over all veh on the entering lane
467                inflow_enter = 0
468                queue_enter = 0
469                for vehid in vehicles_on_entering_lane:
470                    # check which ones are not queueing
471                    if self.sumo_interface.vehicle.getSpeed(vehid) > 0.1:
472                        inflow_enter += 1
473                    # remainder of cars are queuing
474                    else:
475                        queue_enter += 1
476                # save the amount of vehicles that queue to enter the pc and on the seco
477                # key only needs to be created once
478                if 'queue_enter' not in self.traffic_light_data[tl_id] or not self.traff
479                    self.traffic_light_data[tl_id].setdefault('queue_enter', [])
480                    self.traffic_light_data[tl_id].setdefault('inflow_enter', [])
481                self.traffic_light_data[tl_id]['queue_enter'] = queue_enter
482                self.traffic_light_data[tl_id]['inflow_enter'] = inflow_enter
483
484            # east side tl
485            if line_counter == 2:
486                vehicles_on_entering_lane = []
487                vehicles_on_sec_link = []
488                # fill up the dictionary 'self.traffic_light_data[]
489                # dict will contain: {line_counter, lanes_into_perim, vehicles_on_enteri
```

98

```
490                    if 'lanes_into_perim' not in self.traffic_light_data[tl_id] or not self.
491                        self.traffic_light_data[tl_id].setdefault('lanes_into_perim', [])
492                    lanes_into_perim = [sorted_lanes[0], sorted_lanes[1]]
493                    self.traffic_light_data[tl_id]['lanes_into_perim'].extend(lanes_into_per

495                    # Loop through each lane in lanes_into_perim and get the vehicles that e
496                    for lane in lanes_into_perim:
497                        if len(self.sumo_interface.lane.getLastStepVehicleIDs(lane)) > 0:
498                            vehicles_on_entering_lane.extend(self.sumo_interface.lane.getLas

500                    # save the vehicle ids of the vehicles on entering and sec lanes
501                    # key only needs to be created once
502                    if 'vehicles_on_entering_lane' not in self.traffic_light_data[tl_id] or
503                        self.traffic_light_data[tl_id].setdefault('vehicles_on_entering_lane
504                    self.traffic_light_data[tl_id]['vehicles_on_entering_lane'].extend(vehic

506                    # loop over all veh on the entering lane
507                    inflow_enter = 0
508                    queue_enter = 0
509                    for vehid in vehicles_on_entering_lane:
510                        # check which ones are not queueing
511                        if self.sumo_interface.vehicle.getSpeed(vehid) > 0.1:
512                            inflow_enter += 1
513                        # remainder of cars are queuing
514                        else:
515                            queue_enter += 1
516                    # save the amount of vehicles that queue to enter the pc and on the seco
517                    # key only needs to be created once
518                    if 'queue_enter' not in self.traffic_light_data[tl_id] or not self.traff
519                        self.traffic_light_data[tl_id].setdefault('queue_enter', [])
520                        self.traffic_light_data[tl_id].setdefault('inflow_enter', [])
521                    self.traffic_light_data[tl_id]['queue_enter'] = queue_enter
522                    self.traffic_light_data[tl_id]['inflow_enter'] = inflow_enter

524            #south side tl
525            if line_counter == 3:
526                    vehicles_on_entering_lane = []
527                    vehicles_on_sec_link = []
528                    # fill up the dictionary 'self.traffic_light_data[]
529                    # dict will contain: {line_counter, lanes_into_perim, vehicles_on_enteri
530                    if 'lanes_into_perim' not in self.traffic_light_data[tl_id] or not self.
531                        # print('self.traffic_light_data[tl_id]' , self.traffic_light_data[t
532                        self.traffic_light_data[tl_id].setdefault('lanes_into_perim', [])
533                    lanes_into_perim = [sorted_lanes[2], sorted_lanes[3]]
534                    self.traffic_light_data[tl_id]['lanes_into_perim'].extend(lanes_into_per

536                    # Loop through each lane in lanes_into_perim and get the vehicles that e
537                    for lane in lanes_into_perim:
538                        if len(self.sumo_interface.lane.getLastStepVehicleIDs(lane)) > 0:
539                            vehicles_on_entering_lane.extend(self.sumo_interface.lane.getLas
```

```
540
541                             # save the vehicle ids of the vehicles on entering and sec lanes
542                             # key only needs to be created once
543                             if 'vehicles_on_entering_lane' not in self.traffic_light_data[tl_id] or
544                                 self.traffic_light_data[tl_id].setdefault('vehicles_on_entering_lane
545                             self.traffic_light_data[tl_id]['vehicles_on_entering_lane'].extend(vehic
546
547                             # loop over all veh on the entering lane
548                             inflow_enter = 0
549                             queue_enter = 0
550                             for vehid in vehicles_on_entering_lane:
551                                 # check which ones are not queueing
552                                 if self.sumo_interface.vehicle.getSpeed(vehid) > 0.1:
553                                     inflow_enter += 1
554                                 # remainder of cars are queuing
555                                 else:
556                                     queue_enter += 1
557
558                             # save the amount of vehicles that queue to enter the pc and on the seco
559                             # key only needs to be created once
560                             if 'queue_enter' not in self.traffic_light_data[tl_id] or not self.traff
561                                 self.traffic_light_data[tl_id].setdefault('queue_enter', [])
562                                 self.traffic_light_data[tl_id].setdefault('inflow_enter', [])
563                             self.traffic_light_data[tl_id]['queue_enter'] = queue_enter
564                             self.traffic_light_data[tl_id]['inflow_enter'] = inflow_enter
565
566                 # avg speed of all vehicles over 90 sec (one step)
567                 self.episode_90_sec_mean.append(np.mean(self.avg_speed_per_sec))
568
569                 print('******** 90 sec mean speed is *********', self.episode_90_sec_mean[-1])
570
571                 # prevent NaN is needed
572                 if len(self.mean_waiting_per_sec) == 0:
573                     self.mean_waiting_time_90_sec = [0]
574                 else:
575                     self.mean_waiting_time_90_sec.append(np.mean(self.mean_waiting_per_sec))
576
577                 # change in waiting time is mean waiting time over 90 sec minus waiting time on prev
578                 self.change_in_waiting_time.append(np.mean(self.mean_waiting_per_sec) - self.prev_wa
579                 self.prev_waiting_time = (np.mean(self.mean_waiting_per_sec))
580
581                 # store the avg speed inside and outside of the perim per step (90 sec)
582                 if len(self.avg_speed_inside_perim_per_sec)>0:
583                     self.avg_speed_inside_perim.append(np.mean(self.avg_speed_inside_perim_per_sec))
584                 if len(self.avg_speed_outside_perim_per_sec) >0:
585                     self.avg_speed_outside_perim.append(np.mean(self.avg_speed_outside_perim_per_sec
586
587                 # clear so that the list can be refilled each step
588                 self.avg_speed_per_sec.clear()
589                 self.mean_waiting_per_sec.clear()
```

```
590            self.avg_speed_inside_perim_per_sec.clear()
591            self.avg_speed_outside_perim_per_sec.clear()
592
593            state = self.get_state()
594            reward = self.compute_reward()
595
596            # make sure not to overwrite the done given by the empty sim check
597            if done == False:
598                done = self.current_step >= self.max_steps
599            # if the network is empty the simulation stops
600
601            end_time_90_sec_loop = time.time()
602            info = {}
603
604            truncated = False
605            if done:
606                # save the avg speed per episode
607                self.avg_speed = np.mean(self.episode_90_sec_mean)
608
609                # Log avg speed to file
610                with open(self.log_file_path, 'a') as f:
611                    f.write(f"{self.episode_counter},{np.mean(self.avg_speed)}\n")
612                print(' avg speed episode is', self.avg_speed, flush=True)
613                self.avg_speeds_per_step.clear()
614                self.episode_90_sec_mean.clear()
615
616                self.end_time = time.time()
617                print('episode time is' , self.end_time - self.start_time)
618
619            # Check for NaNs in the observation
620            for i, value in enumerate(state):
621                if np.isnan(value):
622                    print(f"NaN found in observation at index {i}: {value}")
623            return state, reward, done, truncated, info
624
625        def reorder_lane_pairs(self, lanes):
626            """Reorder each pair of lanes such that lanes ending in '_0' come before '_1'."""
627            reordered_lanes = []
628            for i in range(0, len(lanes), 2):
629                lane_0 = lanes[i]
630                lane_1 = lanes[i+1]
631
632                # Check the numeric suffixes and reorder
633                if '_0' in lane_0:
634                    reordered_lanes.append(lane_0)
635                    reordered_lanes.append(lane_1)
636                else:
637                    reordered_lanes.append(lane_1)
638                    reordered_lanes.append(lane_0)
639
```

```python
640            return reordered_lanes
641
642       # get controlled lanes takes the lanes in a random order which is not useful
643       # the following couple functions use the name indexing of the lanes to order the control
644       # they are ordered clockwise starting at east with 0 index followed by 1 index
645       # 0 is the most right lane (from direction of traffci) and 1 is the lane to the left of
646       def get_direction_suffix(self, lane_name):
647           """Extract the directional suffix from a lane name (U, D, L, R)."""
648           if 'U' in lane_name:
649               return 'U'  # Upper part on the X-plane
650           elif 'D' in lane_name:
651               return 'D'  # Lower part on the X-plane
652           elif 'L' in lane_name:
653               return 'L'  # Left part on the Y-plane
654           elif 'R' in lane_name:
655               return 'R'  # Right part on the Y-plane
656           return None
657
658
659       def extract_coordinates(self, lane_name):
660           """Extract the X and Y coordinates from the lane name."""
661           # Split by 'Y' to separate the X and Y values.
662           # Example: 'X1Y2.5' -> X=1, Y=2.5
663           parts = lane_name.split('Y')
664           x_coord = float(re.findall(r"[-+]?\d*\.\d+|\d+", parts[0][1:])[0])
    # Extract number from 'X' part
665           y_part = re.findall(r"[-+]?\d*\.\d+|\d+", parts[1])[0]   # Extract number before dire
666           y_coord = float(y_part)
667           return x_coord, y_coord
668
669       def clockwise_sort_key(self, lane_name):
670           """Generate a sort key based on direction and coordinates."""
671           # Clockwise order starting from Up ('U') -> Right ('R') -> Down ('D') -> Left ('L')
672           direction_order = {'U': 0, 'R': 1, 'D': 2, 'L': 3}
673           direction = self.get_direction_suffix(lane_name)
674
675           # Sort by direction first, then by coordinates if the direction is the same
676           x, y = self.extract_coordinates(lane_name)
677           return (direction_order[direction], x, y)
678
679       def sort_lanes_clockwise(self, lanes):
680           """Sort lanes clockwise starting from east and reorder lane pairs."""
681           # First sort the lanes clockwise
682           sorted_lanes = sorted(lanes, key=self.clockwise_sort_key)
683
684           # Then reorder the lanes to ensure '_0' comes before '_1'
685           final_lanes = self.reorder_lane_pairs(sorted_lanes)
686
687           return final_lanes
688
```

```
689
690
691     def get_state ( self ):
692         """
693         get_state function that returns the current state of the simulation step
694         State consists of:
695         -avg speed of the complete network
696         -agv speed outside of the perimeter
697         -avg speed inside the perimeter
698         -total waiting time outside
699         -change in waiting time
700
701         Returns :
702         - state: state of the environment as a np.float32 array consisting of the above ment
703         """
704
705         # the previously mentioned values are not yet filled in the step function
706         # when resetting this loop is executed once
707         if self . current_step == 0:
708             total_vehicles_inside = 0
709             total_vehicles_outside = 0
710             mean_waiting_time_network = 0
711
712             car_id_inside_perim = []
713             car_id_outside_perim = []
714             speed_inside_perim = []
715             speed_outside_perim = []
716             vehicle_ids = self . sumo_interface . vehicle . getIDList ()
717
718
719             # Calculate the total speed and number of vehicles inside and outside perimeter
720             for vehicle_id in vehicle_ids :
721                 mean_waiting_time_network += self . sumo_interface . vehicle . getWaitingTime ( vehi
722                 if self . is_vehicle_inside_perimeter ( vehicle_id ):
723                     car_id_inside_perim . append ( vehicle_id )
724                     speed_inside_perim . append ( self . sumo_interface . vehicle . getSpeed ( vehicle_i
725                     total_vehicles_inside +=1
726                 else :
727                     car_id_outside_perim . append ( vehicle_id )
728                     speed_outside_perim . append ( self . sumo_interface . vehicle . getSpeed ( vehicle_
729                     total_vehicles_outside += 1
730
731             # Calculate the average speed
732             avg_speed_inside_perim = np . mean ( speed_inside_perim ) if total_vehicles_inside >
733             avg_speed_outside_perim = np . mean ( speed_outside_perim ) if total_vehicles_outside
734             avg_speed_network = ( np . mean ( speed_inside_perim )+ np . mean ( speed_outside_perim )/2)
735             change_in_waiting_time = 0 #no change as there is no previous step
736         else :
737             # to get more accurate values , the calculations are done in the step function
738             avg_speed_network = self . episode_90_sec_mean [-1]
```

103

```
739              avg_speed_inside_perim = self.avg_speed_inside_perim[-1]
740              avg_speed_outside_perim = self.avg_speed_outside_perim[-1]
741              mean_waiting_time_network = self.mean_waiting_time_90_sec[-1]
742              change_in_waiting_time = self.change_in_waiting_time[-1]
743
744          state = np.array([avg_speed_network, avg_speed_outside_perim, avg_speed_inside_perim
745
746          # Replace NaN values with 0, on rare ocasions one of these is empty due to a car bei
747          # so simulation will go on but some division by 0 errors can occur
748          state = np.nan_to_num(state, nan=0.0)
749
750          return state
751
752      def init_correct_phase_states(self, phases):
753          """
754          Modifies the states of all traffic light phases by changing all 'g' (green with conf
755          We only want G (green without confliction)
756          Parameters:
757          - phases (list of self.sumo_interface.trafficlight.Phase): A list of traffic light p
758
759          Returns:
760          - list of self.sumo_interface.trafficlight.Phase: The modified list of traffic light
761          """
762
763          # all conflicted greens should be removed already in netedit
764          corrected_phases = []
765          for phase in phases:
766              # Modify the phase state to change all 'g' to 'r'
767              new_state = phase.state.replace('g', 'r')
768              # Create a new phase with the modified state and original duration
769              corrected_phase = self.sumo_interface.trafficlight.Phase(phase.duration, new_sta
770              corrected_phases.append(corrected_phase)
771          return corrected_phases
772
773      def apply_action(self, action):
774
775          """
776          Modifies the action order to allign correctly for tl on the perimeter according to t
777          Then forwards the action to self.set_red_green_phases to implement the new tl config
778
779          First action value is always inflow/outflow direction of PN forward+right turn
780          Second action value is always inflow/outflow direction of PN left turn
781          Third action is perpendicular forward+right turn
782          Fourth action is perpendicular left turn
783          Since the action is given to all tl lights on the perim and the phases always start
784          the order must be changed depending on the position of the tl on the perim
785
786          :param action: action given by rl algo
787          """
788          # Change traffic light configurations based on action
```

```
789            traffic_lights = self.sumo_interface.trafficlight.getIDList()
790
791        for tl in traffic_lights:
792            #checks if the tl is on perimeter, else no modification, SUMO preset phases are
793            # line counter is used to determine which line passed the " True" value for the
794            on_line, line_counter = self.is_traffic_light_on_perimeter(tl)
795            # only modify tl which are on perim line
796            if on_line is True:
797                #west side tl
798                if line_counter == 0:
799                    # Create a copy of the action array for the modified action
800                    modified_action = action.copy()
801                    # Swap the first and third values
802                    modified_action[0], modified_action[2] = action[2], action[0]
803                    # Swap the second and fourth values
804                    modified_action[1], modified_action[3] = action[3], action[1]
805                    # print('modified action west side tl', modified_action)
806                    self.set_red_green_phases(modified_action, tl)
807                # north side tl
808                if line_counter ==1:
809                    # already correct order
810                    self.set_red_green_phases(action, tl)
811                # east side tl
812                if line_counter ==2:
813                    # Create a copy of the action array for the modified action
814                    modified_action = action.copy()
815                    # Swap the first and third values
816                    modified_action[0], modified_action[2] = action[2], action[0]
817                    # Swap the second and fourth values
818                    modified_action[1], modified_action[3] = action[3], action[1]
819                    self.set_red_green_phases(modified_action, tl)
820                # south side tl
821                if line_counter ==3:
822                    # already correct
823                    self.set_red_green_phases(action, tl)
824
825    def set_red_green_phases(self, action, tl_id):
826        """
827        Applies the action given by the rl algo for all tl on the perimeter (this check is a
828
829        First action value is always inflow/outflow direction of protected network (PN) forw
830        Second action value is always inflow/outflow direction of PN left turn
831        Third action is perpendicular forward+right turn
832        Fourth action is perpendicular left turn
833        Since the action is given to all tl lights on the perim and the phases always start
834        the order must be changed depending on the position of the tl on the perim
835
836        :param action: action given by rl algo
837        """
838
```

105

```
839            # keep track of the original program which will be modified so that we prevent modif
840        if not hasattr(self, 'original_program'):
841            # Get the current (original) program logic
842            current_program_tuple = self.sumo_interface.trafficlight.getAllProgramLogics(tl_
843            #print('Original TL logic:', current_program_tuple[0])
844
845            current_program = current_program_tuple[0]
846
847            # Correct the phase states
848            corrected_phases = self.init_correct_phase_states(current_program.phases)
849
850            # Create a new traffic light program with the corrected phases
851            corrected_program = self.sumo_interface.trafficlight.Logic(
852                current_program.programID,
853                current_program.type,
854                current_program.subParameter,
855                corrected_phases
856            )
857
858            # Store the modified program
859            # modified is removing the conflicting green phases, action is NOT yet applied
860            self.original_program = corrected_program
861        else:
862            # Use the stored original program for modification
863            current_program = self.original_program
864
865        action_counter = 0
866        modified_phases = []
867        for phase in current_program.phases:
868            # If the phase contains yellow lights ('y'), skip modifying it
869            if 'y' in phase.state or 'Y' in phase.state:
870                modified_phases.append(phase)  # Keep the yellow phase unchanged
871                continue
872            # Modify the duration of the non yellow phases
873            if action_counter < len(action):
874                new_duration = phase.duration * action[action_counter]
875                # prevent durations being smaller than the min green time (should not be pos
876                if new_duration <= self.min_green_time:
877                    new_duration = self.min_green_time
878                action_counter += 1
879            else:
880                new_duration = phase.duration
881            # Create a new phase with the modified duration and the same state
882            modified_phases.append(self.sumo_interface.trafficlight.Phase(new_duration, phas
883
884
885        # Create a new traffic light program with the modified phases
886        program = self.sumo_interface.trafficlight.Logic("0", 0, 0, modified_phases)
887
888        # Set the new traffic light program
```

```
889            self.sumo_interface.trafficlight.setProgramLogic(tl_id, program)
890
891        def is_traffic_light_on_perimeter(self, tl_id):
892            """
893            Checks if a traffic light is on the perimeter
894            :param tl_id: The ID of the traffic light.
895            :return: True if the tl is on the perimeter, False otherwise.
896
897            Note that this approach does not work for non-straight roads
898
899            """
900            lane_coordinates = []
901            x_coordinates = []
902            y_coordinates = []
903
904            # Get a list of lanes controlled by the tl
905            controlled_lanes = self.sumo_interface.trafficlight.getControlledLanes(tl_id)
906
907            # Loop over all lanes and get their start and end coordinates. (all lanes are straig
908            for lane in controlled_lanes:
909                lane_coordinates.append(self.sumo_interface.lane.getShape(lane))
910
911            # Get the x and y coordinates of the end points of the lane (where the tl is),
912            for x in range(len(controlled_lanes)):
913                x_coordinates.append(lane_coordinates[x][-1][0])
914                y_coordinates.append(lane_coordinates[x][-1][1])
915
916            x_coordinate_tl = np.mean(x_coordinates)
917            y_coordinate_tl = np.mean(y_coordinates)
918
919
920            xmin, ymin, xmax, ymax = self.perim_box
921            coordinates_tl = [x_coordinate_tl, y_coordinate_tl]
922
923            # Define the corners of the perimeter
924            corners = [
925                (xmin, ymin),   # Bottom-left corner
926                (xmin, ymax),   # Top-left corner
927                (xmax, ymax),   # Top-right corner
928                (xmax, ymin)    # Bottom-right corner
929            ]
930
931            # Check if the TL is located at any of the corners, corners will not be considered a
932            for corner in corners:
933                if np.isclose(coordinates_tl[0], corner[0], atol=10) and np.isclose(coordinates_
934                    # print('tl id at corner is', tl_id)
935                    return False, None  # TL is at the corner, return False
936
937            # Define the sides of the perimeter box as LineString objects
938            lines = [
```

```
939                LineString([(xmin, ymin), (xmin, ymax)]),  # Left side
940                LineString([(xmin, ymax), (xmax, ymax)]),  # Top side
941                LineString([(xmax, ymax), (xmax, ymin)]),  # Right side
942                LineString([(xmax, ymin), (xmin, ymin)])   # Bottom side
943            ]
944
945
946            # Depending on which iteration of line the statement is true, the position of the tl
947            # Each loop of the line in lines correspond to the next lines above so left, top, ri
948            # Check if the point is near any of the lines
949            line_counter = 0
950            for line in lines:
951                if self.point_near_line(coordinates_tl, line):
952                    on_line = True
953                    return on_line, line_counter
954                else:
955                    on_line = False
956                    line_counter += 1
957            # (works)
958            return on_line, line_counter
959
960        def is_vehicle_inside_perimeter(self, vehicle_id):
961            """
962            Check if a vehicle is inside the perimeter.
963
964            :param vehicle_id: The ID of the vehicle.
965            :return: True if the vehicle is inside the perimeter, False otherwise.
966            """
967            # x, y coordinate of vehicle
968            x, y = self.sumo_interface.vehicle.getPosition(vehicle_id)
969            # coordinates of the perimeter
970            xmin, ymin, xmax, ymax = self.perim_box
971            # check if a vehicle is within the perimeter
972            return xmin <= x <= xmax and ymin <= y <= ymax
973
974        def point_near_line(self, point, line, tolerance=20.0):
975            """
976            Checks if the point is within a certain tolerance of a line.
977
978            :param point: A tuple (x, y) representing the point.
979            :param line: A LineString object representing the line.
980            :param tolerance: The tolerance within which the point should be from the line (in m
981            :return: True if the point is within the tolerance of the line, False otherwise.
982            """
983            point_obj = Point(point)
984            # determines the euclidian distance from the point to the line
985            distance = point_obj.distance(line)
986            # If the distance is more than the tolerance it will return False, else True
987            return distance <= tolerance
988
```

```
989     def get_inflow_and_queue_enter_and_vehSec_and_queueSec (self):
990         """
991         Retrieves the amount of vehicles that are driving on a lane that enters the PN
992         return: inflow_enter, queue enter, total_vehicles_in_sec_link, queu_sec (int, int, i
993         """
994         inflow_enter = 0
995         queue_enter = 0
996
997         for tl_id, data in self.traffic_light_data.items():
998             inflow_enter += data.get('inflow_enter')
999             queue_enter += data.get('queue_enter')
1000
1001         return inflow_enter, queue_enter
1002
1003     def normalize(self, value, min_val, max_val):
1004         """
1005         Normalizes the value given a min and max.
1006
1007         :param value: desired value to normalize.
1008         :param min_value: min value.
1009         :param max_val: maximum value.
1010
1011         :return: Normalized value
1012         """
1013         return (value - min_val) / (max_val - min_val)
1014
1015     def compute_reward(self):
1016         """
1017         Computes the reward for a point in time in the simulation. This reward is then norma
1018
1019         :return: Normalized reward
1020         """
1021
1022         # Calculate trip completion rate
1023         # if self.total_departed_vehicles == 0:
1024         #     trip_completion_rate = 0.0  # Avoid division by zero if no vehicles have start
1025         # else:
1026         #     # 4000 comes from generate_trips.py " trip_count"
1027         #     trip_completion_rate = (self.total_arrived_vehicles / 4000) * 100
1028
1029         # # two methods for trip completion rate:
1030         # # total arrived/ total departed so far
1031         # # toal arrived/total trips over whole sim
1032
1033         # inflow_enter, queue_enter = self.get_inflow_and_queue_enter_and_vehSec_and_queueSe
1034         # vehicles_enter = inflow_enter + queue_enter
1035
1036         # x1 = 0.9
1037         # x2 = -0.1
1038         # x3 = -0.1
```

```
1039
1040        # # get the most recent value for the avg speed over 90 sec to find the avg speed ov
1041
1042        # change_in_waiting_time = self.change_in_waiting_time[-1]
1043
1044        # print(' TCR', trip_completion_rate)
1045        # # use tanh transform to make the value between -1 and 1. Shift the range to be bet
1046        # normalized_trip_completion_rate = (np.tanh(trip_completion_rate)+1)/2
1047        # normalized_vehicle_enter = (np.tanh(vehicles_enter)+1)/2
1048
1049        # calculate reward
1050        # reward = x1*normalized_trip_completion_rate + x2*normalized_vehicle_enter + x3
1051
1052        avg_speed = self.episode_90_sec_mean[-1]
1053        reward = avg_speed
1054        return reward
1055
1056        # #basic reward
1057        # reward = avg_speed
1058        # self.reward_list.append(reward)
1059        # #print(' reward list is', self.reward_list)
1060
1061        # return reward
1062
1063        # min_speed = 0
1064        # max_speed = 8
1065
1066        # min_waiting_time = 0
1067        # max_waiting_time = 3000
1068
1069        # normalized_avg_speed = self.normalize(avg_speed, min_speed, max_speed )
1070        # normalized_waiting_time =self.normalize(change_in_waiting_time, min_waiting_time,
1071        # use normalized avg speed and waiting time in future
1072
1073        # uncommend for weighted normalized reward
1074        # Raw reward calculation
1075
1076        # reward = avg_speed #- 0.1 * change_in_waiting_time
1077        # reward  = avg_speed #- 0.1*normalized_waiting_time
1078        # return reward
1079        # # Update running mean and variance for normalization
1080        # self.num_rewards += 1
1081        # old_mean = self.reward_mean
1082        # self.reward_mean += (reward - self.reward_mean) / self.num_rewards
1083        # self.reward_var += (reward - old_mean) * (reward - self.reward_mean)
1084
1085        # reward_std = np.sqrt(self.reward_var / (self.num_rewards - 1)) if self.num_rewards
1086
1087        # # Normalize the reward using running mean and standard deviation
1088        # normalized_reward = -((reward - self.reward_mean) / (reward_std + 1e-8))
```

```
1089
1090            # return normalized_reward
1091
1092        def close(self):
1093            """ closes the self.sumo_interface connection """
1094            # close sumo connection
1095            self.sumo_interface.close()
1096
1097        def plot_avg_speed_per_episode(self):
1098            """Plot the average speed per episode."""
1099            plt.figure(figsize=(10, 6))
1100            plt.plot(self.avg_speed_per_episode, marker='o', linestyle='-', color='b')
1101            plt.title('Average Speed per Episode')
1102            plt.xlabel('Episode')
1103            plt.ylabel('Average Speed')
1104            plt.grid(True)
1105            plt.show()
1106
1107
1108    class SumoGymEnv(gymnasium.Env):
1109        """ This class does the RL part """
1110
1111        def __init__(self, sumo_cfg, max_steps, rank):
1112            """"initializes the class"""
1113            super(SumoGymEnv, self).__init__()
1114            #print(' action space has been modified for error testing')
1115            #lower_limit_action = np.array([1, 1, 1, 1])
1116            #upper_limit_action = np.array([1, 1, 1, 1])
1117            lower_limit_action = np.array([0.5, 0.75, 0.5, 0.75])
1118            upper_limit_action = np.array([1.5, 1.8, 1.5, 1.8])
1119            self.sumo_env = SumoEnvironment(sumo_cfg, max_steps, rank)
1120            self.action_space = gymnasium.spaces.Box(low=lower_limit_action, high=upper_limit_ac
1121            self.observation_space = gymnasium.spaces.Box(low=0, high=np.inf, shape=(5,), dtype=
1122
1123
1124        def reset(self, *, seed=None, options=None):
1125            """ Reset function, calls the sumoEnvironment reset function"""
1126            # Seed the environment's RNGs if applicable
1127            if seed is not None:
1128                super().reset(seed=seed)
1129                np.random.seed(seed)
1130            return self.sumo_env.reset()
1131
1132        def step(self, action):
1133            """ Step function, calls the SumoEnvironment step function"""
1134            return self.sumo_env.step(action)
1135
1136        def close(self):
1137            """ close function, calls the SumoEnvironment close function"""
1138            self.sumo_env.close()
```

```
1139
1140
1141    # class CustomTensorboardCallback(BaseCallback):
1142    #       """
1143    #       Custom callback for logging additional metrics to TensorBoard, such as policy loss, va
1144    #       """
1145    #       def __init__(self, verbose=0):
1146    #           super(CustomTensorboardCallback, self).__init__(verbose)
1147
1148    #       def _on_step(self) -> bool:
1149    #           print(' logging test print')
1150    #           print(' self.model.policy is', self.model.policy)
1151    #           # Access the model's rollout buffer to get the advantages, rewards, and actions
1152    #           rollout_buffer = self.model.rollout_buffer
1153
1154    #           # Mean advantage, policy loss, and value loss are calculated during training
1155    #           mean_advantage = rollout_buffer.advantages.mean().item()
1156    #           mean_action = rollout_buffer.actions.mean().item()
1157    #           # Retrieve the policy and value losses from the PPO training process
1158    #           # Accessing the PPO losses requires some manual computation during the optimizatio
1159    #           # Get the distribution (actor network output)
1160    #           #actions_dist = self.model.policy.dist
1161    #           if hasattr(self.model, 'policy'):
1162    #               policy_loss = self.model.policy.loss.mean().item() if hasattr(self.model.polic
1163    #               value_loss = self.model.policy.value_loss.mean().item() if hasattr(self.model.
1164
1165    #           else:
1166    #               policy_loss = None
1167    #               value_loss = None
1168    #               mu, sigma = None, None
1169
1170    #           # Log the statistics to TensorBoard
1171    #           if policy_loss is not None and value_loss is not None:
1172    #               self.logger.record("ppo/policy_loss", policy_loss)
1173    #               self.logger.record("ppo/value_loss", value_loss)
1174
1175    #           self.logger.record("ppo/mean_advantage", mean_advantage)
1176    #           self.logger.record("ppo/mean_action", mean_action)
1177    #           print('policy loss sent to tensorboard is', policy_loss, flush=True)
1178    #           print('value loss sent to tensorboard is', value_loss, flush=True)
1179    #           print('mean advantage sent to tensorboard is', mean_advantage, flush=True)
1180    #           print('mean action sent to tensorboard is', mean_action, flush=True)
1181    #           print('final mu sent to tensorboard is', mu[-1], flush=True)
1182    #           print('final sigma sent to tensorboard is', sigma[-1], flush=True)
1183
1184    #           return True
1185
1186
1187
1188    if __name__ == '__main__':
```

```
1189
1190     # Initialize SUMO gym environment
1191     sumo_cfg = "sumo_rl/nets/PC-Alex/city_network_main.sumocfg"
1192     max_steps = 3000     #max steps, warm up time is already considered in reset, multiple of
1193     # cars are spawned until 3600, the remainder is to empty the network
1194     print(' max steps has been changed to use dit4 reward')
1195     # env = SumoGymEnv(sumo_cfg, max_steps)
1196
1197     # Number of parallel environments
1198     num_envs = 1
1199     print(' number of envs changed back to 1 for error testing')
1200
1201     # Create the vectorized environment
1202     env_fns = [SumoEnvironment.make_sumo_env(sumo_cfg, max_steps, rank=i) for i in range(num
1203
1204     # creates vectorized env
1205     env = SubprocVecEnv(env_fns)
1206
1207     #experiment settings
1208     max_iters = 50
1209     TIMESTEPS = max_steps
1210
1211     # logging settings
1212     log_name = f"PPO_city_simple_iters_{max_iters}_timesteps_{TIMESTEPS}"
1213
1214     models_dir = f"models/{log_name}/"
1215     logdir = "logs_ppo_city/ppo_city_simple" + time.strftime("%Y-%m-%d_%H-%M-%S")
1216
1217     if not os.path.exists(models_dir):
1218         os.makedirs(models_dir)
1219
1220     if not os.path.exists(logdir):
1221         os.makedirs(logdir)
1222
1223     # modify the neural network
1224     net_arch = [5, 256, 256, 256, 256, 4]
1225
1226
1227     # Initialize the PPO model with your environment
1228     # model = PPO('MlpPolicy', env, batch_size=32, ent_coef=0.1, learning_rate=0.001, n_step
1229     verbose=1, tensorboard_log='logs_ppo_city/ppo_city_simple'+ time.strftime("%Y-%m-%d_%H-%M-%S
1229     model = PPO('MlpPolicy', env, batch_size=64, ent_coef=0.2, learning_rate=0.01, n_steps=1
1229     verbose=1, tensorboard_log=logdir)
1230
1231
1232     # Set the model in one of the environments for logging
1233     env_fns[0]().sumo_env.set_model(model)  # Set the model in the first environment for log
1234
1235     # TensorBoard logging setup
1236
```

```
1237        summary_writer = tf.summary.create_file_writer(logdir)
1238        # Define the callback
1239        # avg_speed_logging_callback = AvgSpeedLoggingCallback(summary_writer=summary_writer)
1240
1241        iters =0
1242
1243        # Training loop
1244        while iters <= max_iters:
1245            # Start timing the episode
1246            start_time = time.time()
1247            iters += 1
1248            # model.learn(total_timesteps=TIMESTEPS, reset_num_timesteps=False, log_interval=1,
1249            model.learn(total_timesteps=TIMESTEPS, reset_num_timesteps=False, log_interval=1)
1250            #model.learn(total_timesteps=TIMESTEPS, reset_num_timesteps=False, tb_log_name=log_n
1251            # End timing the episode
1252            end_time = time.time()
1253
1254            # Calculate and print the duration of the episode
1255            episode_duration = end_time - start_time
1256            print(f"Time taken for episode {iters}: {episode_duration:.2f} seconds")
1257
1258            model.save(f"{models_dir}/{log_name}")
1259
1260
1261
1262 # logging residual code
1263 # # Access the PPO's `rollout_buffer` and get the advantage from the buffer
1264 #                advantages = self.model.rollout_buffer.advantages
1265 #                mean_advantage = advantages.mean().item()  # Log the mean of the advantages
1266
1267 #                # Get policy loss and value loss
1268 #                policy_loss = self.model.policy.loss
1269 #                value_loss = self.model.policy.value_loss
1270
1271 #                # Log the policy (actions chosen by the agent)
1272 #                observations = self.locals['rollout_buffer'].observations
1273 #                actions = self.model.predict(observations, deterministic=False)[0]
     # Get actions from the policy
1274 #                mean_action = actions.mean()  # Optionally, log the mean of the actions taken
1275
1276
1277 #                policy = self.model.policy
1278
1279 #                # Get the distribution (actor network output)
1280 #                actions_dist = policy.dist
1281 #                mu = actions_dist.mean
1282 #                sigma = actions_dist.stddev
1283
1284 #                 # Log to TensorBoard
1285 #                with self.summary_writer.as_default():
```

114

```
1286   #                     tf.summary.scalar('PPO/policy_loss', policy_loss.item(), step=self.current
1287   #                     tf.summary.scalar('PPO/value_loss', value_loss.item(), step=self.current_s
1288   #                     tf.summary.scalar('PPO/mean_advantage', mean_advantage, step=self.current_
1289   #                     tf.summary.scalar('PPO/mean_action', mean_action, step=self.current_step)
1290   #                     # Log the mean and stddev per action value (action array consists of 4 act
1291   #                     for i in range(mu.shape[0]):
1292   #                         tf.summary.scalar(f'PPO/action_mu_{i}', mu[i], step=self.current_step)
1293   #                         tf.summary.scalar(f'PPO/action_sigma_{i}', sigma[i], step=self.current
1294
1295   #                print('policy loss sent to tensorboard is', policy_loss.item(), flush=True)
1296   #                print('value loss sent to tensorboard is', value_loss.item(), flush=True)
1297   #                print('mean advantage sent to tensorboard is', mean_advantage, flush=True)
1298   #                print('mean action sent to tensorboard is', mean_action, flush=True)
1299   #                print('final mu sent to tensorboard is', mu[-1], flush=True)
1300   #                print('final sigma sent to tensorboard is', sigma[-1], flush=True)
```

```
1    from joblib import load, dump
2    import numpy as np
3    import ast
4    import pandas as pd
5    import re
6    import matplotlib.pyplot as plt
7    from sklearn.neural_network import MLPRegressor
8    from sklearn.model_selection import train_test_split
9    from sklearn.preprocessing import StandardScaler
10   from sklearn.preprocessing import RobustScaler
11   from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
12   from sklearn.model_selection import cross_val_score, KFold
13   from sklearn.preprocessing import MinMaxScaler
14   from scipy.stats import gaussian_kde
15   from sklearn.metrics import mean_absolute_error
16
17   "This file is used for fine tuning the NN, used in the iterative supervised learning approac
18
19
20   def load_existing_model(model_filename, x_scaler_filename, y_scaler_filename):
21       """Load an existing model and scalers from saved files."""
22       model = load(model_filename)
23       x_scaler = load(x_scaler_filename)
24       y_scaler = load( y_scaler_filename)
25       return model, x_scaler, y_scaler
26
27   def preprocess_new_data(df, x_scaler, y_scaler):
28       """Preprocess and scale new data to match existing model."""
29       # Apply the same preprocessing steps as used in the original training data
30       def fix_format(data_string):
31           data_string = data_string.strip()
32           data_string = re.sub(r'\s+', ' ', data_string)
33           fixed_string = re.sub(r'(\d)\s+(\d)', r'\1, \2', data_string)
```

```
34              return fixed_string
35
36      df["state"] = df["state"].apply(fix_format).apply(ast.literal_eval)
37      df["policy"] = df["policy"].apply(fix_format).apply(ast.literal_eval)
38
39      def extract_inner_value(value):
40          return float(value[len("exp("):-1])
41      df['traffic_dist'] = df['traffic_dist'].apply(extract_inner_value)
42
43      states = df["state"].tolist()
44      policies = df["policy"].tolist()
45      traffic_dist = np.array(df["traffic_dist"].tolist()).reshape(-1, 1)
46      X_train = np.hstack((states, policies, traffic_dist))
47
48      # mutliple rewards were attempted, this is the best performing one
49      reward_number = 2
50      y_train = np.array([ast.literal_eval(reward)[reward_number] for reward in df["reward"].t
51
52      # Scale the new data using the original scalers
53      X_train_scaled = x_scaler.transform(X_train)
54      y_train_scaled = y_scaler.transform(y_train.reshape(-1, 1)).flatten()
55
56      return X_train_scaled, y_train_scaled
57
58  def update_model_with_new_data(model, X_train, y_train, epochs=5):
59      """Update the model using new data."""
60      for epoch in range(epochs):
61          model.partial_fit(X_train, y_train)  # Use partial_fit to update model with new data
62          train_loss = np.mean((model.predict(X_train) - y_train) ** 2)
63          print(f"Epoch {epoch + 1}/{epochs} - New Data Training Loss: {train_loss:.4f}")
64      return model
65
66
67  """" MODIFY HERE"""
68  # Load the existing model and scalers
69  model, x_scaler, y_scaler = load_existing_model('iter4_nn_model_rew2.joblib', 'x_scaler_rew2
70  # model, x_scaler, y_scaler = load_existing_model('iter4_nn_model_state2.joblib', 'x_scaler_
71
72  # Load and preprocess new data
73  new_data_filename = "sumo_training_data_iter5.csv"
74  new_validation_filename = "sumo_validation_data_iter5.csv"
75
76  # new_data_filename = "sumo_training_data_state1_iter5.csv"
77  # new_validation_filename = "sumo_validation_data_state1_iter5.csv"
78
79
80  """"""""
81
82  new_data_df = pd.read_csv(new_data_filename, names=["state", "policy", "reward", "traffic_di
83  new_validation_data_df = pd.read_csv(new_validation_filename, names=["state", "policy", "rew
```

```
84
85  X_train, y_train = preprocess_new_data(new_data_df, x_scaler, y_scaler)
86  X_valid, y_valid = preprocess_new_data(new_validation_data_df, x_scaler, y_scaler)
87
88  # Update model with new data
89  updated_model = update_model_with_new_data(model, X_train, y_train)
90
91  # save the updated model
92  dump(updated_model, 'iter5_nn_model_rew2.joblib')
93  dump(x_scaler, 'x_scaler_rew2_iter5.joblib')
94  dump(y_scaler, 'y_scaler_rew2_iter5.joblib')
95
96  # dump(updated_model, 'iter5_nn_model_state2.joblib')
97  # dump(x_scaler, 'x_scaler_state2_iter5.joblib')
98  # dump(y_scaler, 'y_scaler_state2_iter5.joblib')
99
100 # Implement cross-validation
101 kf = KFold(n_splits=5, shuffle=True, random_state=42)
102 cross_val_scores = cross_val_score(model, X_train, y_train, cv=kf, scoring='neg_mean_squared
103
104 print("Cross-Validation MSE for each fold:", -cross_val_scores)
105 print("Cross-Validation Mean MSE:", -np.mean(cross_val_scores))
106
107 # Evaluate the model
108 train_score = model.score(X_train, y_train)
109 valid_score = model.score(X_valid, y_valid)
110 print(f"Training score: {train_score}")
111 print(f"Test score: {valid_score}")
112
113 # Inverse transform predictions for interpretation
114 predictions = model.predict(X_valid)
115 original_scale_predictions = y_scaler.inverse_transform(predictions.reshape(-1, 1)).flatten(
116 original_scale_y_valid = y_scaler.inverse_transform(y_valid.reshape(-1, 1)).flatten()
117
118 # Print the first 5 original-scale true values and predicted values
119 print("First 5 true values (original scale):", original_scale_y_valid[:5])
120
121 print("First 5 predicted values (original scale):", original_scale_predictions[:5])
122
123 plt.figure(figsize=(8, 6))
124 plt.scatter(original_scale_y_valid, original_scale_predictions, alpha=0.6)
125 plt.plot([original_scale_y_valid.min(), original_scale_y_valid.max()],
126         [original_scale_y_valid.min(), original_scale_y_valid.max()], 'r--', lw=2)
127 plt.xlabel('True Values')
128 plt.ylabel('Predicted Values')
129 plt.title('True vs. Predicted Values')
130 plt.show()
131
132 mae = mean_absolute_error(original_scale_y_valid, original_scale_predictions)
133 print("Mean Absolute Error (MAE):", mae)
```

```
134
135  rmse = np.sqrt(np.mean((original_scale_y_valid - original_scale_predictions) ** 2))
136  print("Root Mean Squared Error (RMSE):", rmse)
137
138  # Calculate percentage errors
139  percentage_errors = np.abs((original_scale_y_valid - original_scale_predictions) / original_
140
141  # Perform Kernel Density Estimation to get a smooth curve
142  kde = gaussian_kde(percentage_errors, bw_method=0.3)  # Bandwidth can be adjusted for smooth
143  x_range = np.linspace(0, percentage_errors.max(), 100)
144
145  # Plot smooth density line of the percentage errors
146  plt.figure(figsize=(10, 6))
147  plt.plot(x_range, kde(x_range), label='Density of Percentage Error')
148  plt.xlabel("Percentage Error (%)")
149  plt.ylabel("Density")
150  plt.title("Density Plot of Percentage Errors in Validation Predictions")
151  plt.legend()
152  plt.grid(True)
153  plt.show()
154
155  # Calculate absolute errors for each data point
156  errors = np.abs(original_scale_y_valid - original_scale_predictions)
157
158  # Sort errors in descending order
159  sorted_errors = np.sort(errors)[::-1]
160
161  # Calculate cumulative contribution of errors as percentage
162  cumulative_error = np.cumsum(sorted_errors) / np.sum(sorted_errors) * 100
163
164  # Calculate percentage of data points
165  data_points_percentage = np.arange(1, len(errors) + 1) / len(errors) * 100
166
167  # Calculate absolute percentage errors
168  percentage_errors = np.abs((original_scale_y_valid - original_scale_predictions) / original_
169
170  # Define error thresholds
171  thresholds = np.arange(0, 50, 5)  # Adjust range and step as needed (e.g., up to 50% with 5%
172  percentages_below_threshold = [(percentage_errors <= t).mean() * 100 for t in thresholds]
173
174  # Plot the cumulative distribution of errors
175  plt.figure(figsize=(10, 6))
176  plt.plot(thresholds, percentages_below_threshold, marker='o', linestyle='-', color='b')
177  plt.xlabel('Error Threshold (%)')
178  plt.ylabel('Percentage of Predictions Below Threshold (%)')
179  plt.title('Cumulative Distribution of Prediction Errors')
180  plt.grid(True)
181  plt.ylim(0, 100)
182  plt.show()
183
```

```
184  # Calculate absolute percentage errors
185  absolute_percentage_errors = np.abs((original_scale_y_valid - original_scale_predictions) /
186
187  # Calculate the mean absolute percentage error
188  mape = np.mean(absolute_percentage_errors)
189  print("Mean Absolute Percentage Error (MAPE):", mape)
```