

DrPlanner: Diagnosis and Repair of Motion Planners Using Large Language Models

Yuanfei Lin^{1,2}, Chenran Li², Mingyu Ding², Masayoshi Tomizuka², Wei Zhan², and Matthias Althoff¹

Abstract—Motion planners are essential for the safe operation of automated vehicles across various scenarios. However, no motion planning algorithm has achieved perfection in the literature, and improving its performance is often time-consuming and labor-intensive. To tackle the aforementioned issues, we present DrPlanner, the first framework designed to automatically diagnose and repair motion planners using large language models. Initially, we generate a structured description of the planner and its planned trajectories from both natural and programming languages. Leveraging the profound capabilities of large language models in addressing reasoning challenges, our framework returns repaired planners with detailed diagnostic descriptions. Furthermore, the framework advances iteratively with continuous feedback from the evaluation of the repaired outcomes. Our approach is validated using search-based motion planners; experimental results highlight the need of demonstrations in the prompt and the ability of our framework in identifying and rectifying elusive issues effectively.

I. INTRODUCTION

Motion planners for automated vehicles are responsible for computing safe, physically feasible, and comfortable motions [1]. However, to the best of our knowledge, no universal algorithm currently exists that can safely and reliably solve the motion planning problem in all scenarios (see Sec. I-A.1). Therefore, it is crucial to continuously evaluate and enhance the performance of a motion planner during its development. A major challenge is the excessive manual effort required, which entails diagnosing the planner based on a variety of critical test scenarios and evaluation metrics. This process requires not only deep expertise in motion planner functionalities but also a comprehensive understanding of how various aspects of the algorithm correlate with performance. Furthermore, the discrepancy between the design of the algorithm and its practical implementation is another significant factor to consider. To address these challenges, we establish a nuanced framework that leverages the remarkable emergent abilities of large language models (LLMs) [2]–[4]

The authors gratefully acknowledge partial financial support by the German Federal Ministry for Digital and Transport (BMDV) for the project KoSi, by the European Research Council (ERC) for the project justITSELF, under Grant Agreement No. 817629, and by the Berkeley DeepDrive. The work was developed during Y. Lin’s visit to the University of California, Berkeley. (*Corresponding author: W. Zhan.*)

¹ Y. Lin and M. Althoff are with the School of Computation, Information and Technology, Technical University of Munich, 85748 Garching, Germany. {yuanfei.lin, althoff}@tum.de

² Y. Lin, C. Li, M. Ding, M. Tomizuka, and W. Zhan are with the Department of Mechanical Engineering, University of California, Berkeley, CA 94720, USA. {chenran.li, myding, tomizuka, wzhan}@berkeley.edu

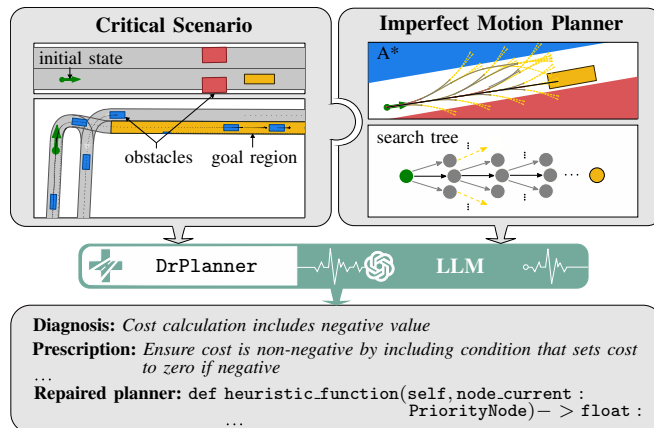


Fig. 1: An example usage of DrPlanner: In a critical scenario, our imperfect motion planner plans a trajectory. The description of the trajectory and the planner is then fed into DrPlanner. By harnessing the strengths of LLMs in understanding common sense and programming languages, we adeptly diagnose and repair the deficiencies within the planner.

to automatically provide and apply diagnostic solutions for a motion planner, which is shown in Fig. 1.

A. Related Work

Subsequently, we review the literature on the imperfections of motion planners, automated software repair, and the application of LLMs in motion planning.

1) *Imperfections of Motion Planners*: Although many motion planning algorithms can tackle a diverse range of tasks, they often face issues related to probabilistic completeness, computational complexity, or real-time constraints in finding the optimal solution [1], [5], [6]. For example, the authors of [7] examine the solvability of the planning problem using multiple trials and a specified time budget. They demonstrate that both their planner and state-of-the-art alternatives [8], [9] might fail to find a feasible solution, even when one exists. Moreover, most previous studies do not benchmark motion planners across a variety of scenarios, even though they are developed with a few use cases in mind. This lack of comprehensive evaluation makes it challenging to assess and compare the algorithms, let alone improve them. Besides, advanced deep learning algorithms [10]–[12] have been applied over the past decades to handle complex scenarios and learn from experience. However, guaranteeing safety, rule compliance, and social compatibility of motion planners remains a challenge [13]–[15].

2) *Automated Software Repair*: With the increasing complexity and size of software, automatic debugging and repair techniques have been developed to reduce the extensive manual effort required to fix faults and to improve quality [16].

For instance, human-designed templates are used to repair certain types of bugs in code [17]–[21], but their effectiveness is limited to the hard-coded patterns. To overcome these limitations, deep-learning-based approaches utilize neural machine translation [22] to learn from existing patches, treating the repaired code as a translation of the buggy one [23]–[26]. However, the performance of these approaches is limited by the quality and quantity of the training data as well as its representation format [27]. As LLMs have shown emergent abilities in solving programming tasks [28]–[31], they are applied for generating program patches [32]–[34], self-debugging [35], [36], and cleaning code [37]. Unlike simply maintaining functional equivalence, we aim to both rectify imperfections and boost the performance of the planning algorithms. The aspect of performance improvement aligns with [38], but our work distinguishes itself by focusing on motion planners with a larger codebase. Another branch of work focuses on repairing the outcome of given software [39]–[41] or addressing specified diagnostic criteria [42].

3) *Language Models for Motion Planning*: With their indispensable role of common sense reasoning and generalization [43]–[45], LLMs have been applied in motion planning for autonomous driving to make high-level decisions [46]–[50], generate driving trajectories [51], [52] or provide control signals directly [53]–[55]. However, the refinement of motion planners themselves is still driven by the nuanced intuition of humans and by real traffic data. In this work, LLMs serve to bridge this gap by emulating human-like problem-solving strategies, offering strategic guidance in analyzing complex motion planners.

B. Contributions

In this work, we introduce DrPlanner, the first framework to autonomously diagnose and repair motion planners, harnessing the power of LLMs that improve as they scale with additional data and model complexity. In particular, our contributions are:

- 1) establishing a structured and modular description for motion planners across both natural and programming language modalities to exploit the capabilities of LLMs for diagnosis and repair;
- 2) leveraging the in-context learning capabilities of LLMs by providing demonstrations to the model at the point where it infers diagnostic results;
- 3) and enhancing the understanding of underlying improvement mechanisms by generating continuous feedback in a closed-loop manner.

The remainder of this work is structured as follows: Sec. II lists necessary preliminaries. The proposed framework for diagnosing and repairing motion planners is described in Sec. III. We demonstrate the benefits of our approach in Sec. IV and conclude the paper in Sec. V.

II. PRELIMINARIES

A. Motion Planning

We refer to the vehicle for which trajectories are planned as the *ego vehicle*. As illustrated in Fig. 2, motion planning

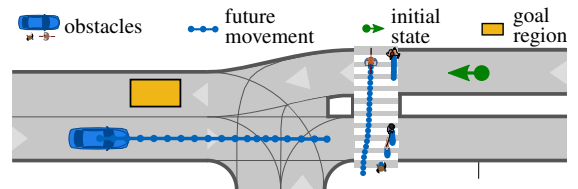


Fig. 2: Exemplary motion planning problem, where the ego vehicle needs to travel from its initial state to reach the goal region safely and efficiently.

algorithms are tasked with ensuring that the ego vehicle travels from an initial state to a goal region within a specified time [56]. Additionally, the solution, denoted by χ , must satisfy common and safety-relevant requirements, such as being drivable, collision-free, and rule-compliant [40], [57]. Meanwhile, the motion planner typically minimizes a given objective function $J(\chi)$, e.g., by penalizing the travel time or passenger discomfort [1, Sec. IV]. Finally, we denote a motion planner by M and a motion planning problem by P .

B. Prompt Engineering for LLMs

The technique of using a textual string ℓ to condition LLMs for probabilistic predictions is referred to as *prompting* [58]. This approach enables LLMs to be pretrained on a massive amount of data and subsequently adapt to new use cases with few or no labeled data. To enhance the in-context learning capabilities, the prompt may include a few human-annotated examples of the task, known as *few-shot prompting* [2], or utilize chain-of-thought reasoning [43], [59]. We divide the input prompt ℓ into two components: the system prompt ℓ_{system} , which outlines the task for the LLMs, and the user prompt ℓ_{user} , providing context for the diagnostic task. The labels, manual inputs, and automatically generated content within the prompt are marked with angle brackets, square brackets, and curly brackets, respectively. The output consists of both a list of diagnosis-prescription pairs and patched programs, collectively denoted by ℓ_{dp} and p_p . It is important to note that LLMs typically have a limit on the number of tokens [60] they can process, which essentially means there is a maximum length for the prompt.

III. DRPLANNER

If an LLM were asked, “*How can the performance of a motion planner be improved?*”, it might respond, “*It involves a combination of software optimization, algorithmic enhancements, and hardware improvements.*” Although this response is a reasonable completion for the prompt, it is not necessarily actionable in diagnosing the downstream planner with specific setups. Therefore, we adopt careful prompt engineering with a nuanced diagnostic description, as explained in this section. We begin by introducing the overall algorithm, followed by its details.

A. Overall Algorithm

A general overview of using DrPlanner to diagnose and repair planners is presented in Fig. 3 and Alg. 1. Before initiating the process, the user must fill in the placeholders enclosed in square brackets with the required manual input.

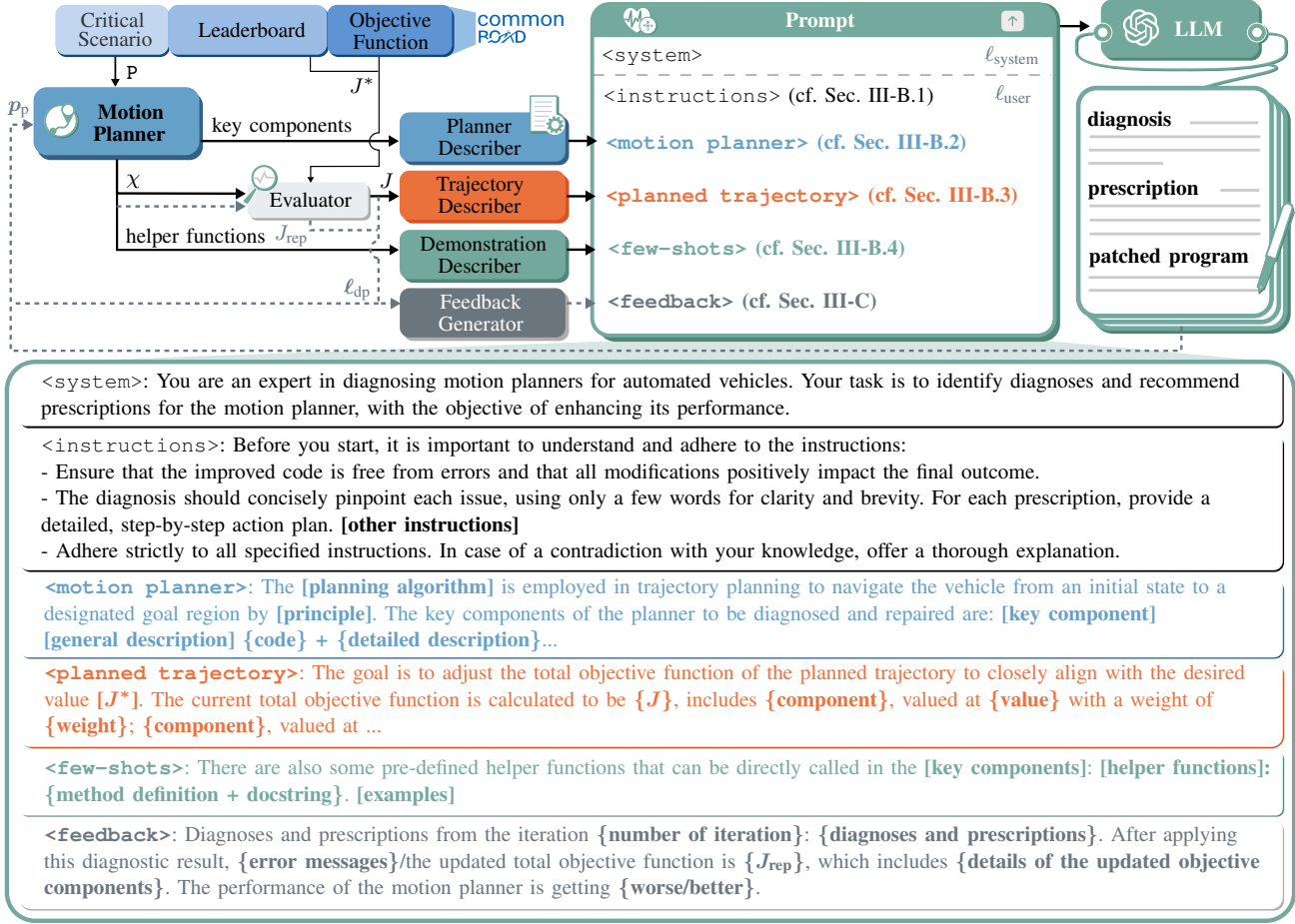


Fig. 3: Overview of the DrPlanner framework. The process starts with obtaining a planned trajectory for the planning problem with the given motion planner. Then, the planned trajectory is evaluated by the objective function. Afterwards, the description for the planner is generated and used to prompt an off-the-shelf LLM to generate the diagnoses and prescriptions for the planner, along with the patched programs. After applying the patches, the evaluation of the updated planner is incorporated back into the prompt as feedback to continuously enhance the diagnostic performance (marked by dashed arrows).

Algorithm 1 DIAGNOSEANDREPAIRPLANNER

Input: planning problem P , motion planner M , target value J^* , system prompt ℓ_{system} , LLM

Output: diagnoses and prescriptions ℓ_{dp}^* , repaired planner M_{rep}^*

- 1: $\chi \leftarrow M.PLAN(P)$
- 2: $J \leftarrow EVALUATE(\chi)$
- 3: $\ell_{user} \leftarrow DESCRIBE(M, J, J^*)$ ▷ Sec. III-B
- 4: $J_{min} \leftarrow J, \ell_{dp}^* \leftarrow \emptyset, M_{rep}^* \leftarrow \emptyset$
- 5: **while** not REACHTOKENLIMIT(LLM) and $J_{min} - J^* > \epsilon$ **do**
- 6: $(\ell_{dp}, p_p) \leftarrow LLM.QUERY(\ell_{system}, \ell_{user})$ ▷ Sec. III-C
- 7: $M_{rep} \leftarrow REPAIR(M, p_p)$
- 8: $\chi \leftarrow M_{rep}.PLAN(P)$
- 9: $J_{rep} \leftarrow EVALUATE(\chi)$
- 10: $\ell_{user} \leftarrow ADDFEEDBACK(\ell_{user}, J, J_{rep}, \ell_{dp})$ ▷ Sec. III-C
- 11: **if** $J_{rep} < J_{min}$ **then**
- 12: $J_{min} \leftarrow J_{rep}, \ell_{dp}^* \leftarrow \ell_{dp}, M_{rep}^* \leftarrow M_{rep}$
- 13: **end if**
- 14: **end while**
- 15: **return** ℓ_{dp}^*, M_{rep}^*

For a given scenario, the motion planner M is first deployed to address the associated planning problem P (see line 1). Subsequently, the planned trajectory χ is evaluated against the objective function J (see line 2). Following this, a diagnostic description ℓ_{user} encompassing the diagnostic instructions,

the description of the planner, the evaluation of the trajectory, and the few-shot examples are formulated (see line 3). This description, along with the system prompt ℓ_{system} , is then fed into the LLM (see line 6). The structure of the input prompt is illustrated in the center of the framework in Fig. 3. Afterwards, the obtained patched programs are applied to the motion planner by integrating the modifications into the existing codebase (see line 7). However, it is important to note that the output generated may include errors such as hallucinations and inaccurate analyses [61]. To mitigate these issues, we employ an iterative prompting strategy repeatedly refining the process. The iteration is terminated when a notable improvement in the planner is observed, e.g., when the difference between J_{min} and a target value J^* is smaller than a threshold $\epsilon \in \mathbb{R}_+$, or when the token limit of the LLM is reached (see lines 5-14). Finally, the repaired planner demonstrating the best improvement, if any, along with the corresponding diagnoses and prescriptions, is returned (see line 15). Another regime is to finetune the LLM to the diagnostic and repair task. However, to date, there exists no dataset containing input-output examples as the cumbersome improvement process of motion planners is typically neither open-sourced nor well documented. Additionally, finetuning

usually only provides modest improvement in solving challenging and complex tasks compared to in-context learning [36], [37], [59].

B. Diagnostic Description

As discussed in Sec. II-B, prompt design is challenging, particularly when considering the limited information about the diagnostic object in the pretrained LLM. To enhance reasoning outcomes, we design a structured and comprehensive description of the motion planner, emulating the process of a real doctor. Its overall skeleton is depicted in the lower part of Fig. 3. As we assume that the motion planner internally handles goal-reaching and drivability-checking of the trajectory in the scenario (cf. Sec. II-A), a detailed description of the scenario, motion planning problem, and trajectory is omitted in the prompt. Alternatively, these tasks can be addressed by additional modules, such as those employing LLM-embedded agents (cf. Sec. I-A.3).

1) *Instructions*: The instruction provides general guidance for the LLM, detailing the expected output and reasoning constraints. In addition, we can include the commonly used rule-of-thumb from expert knowledge. For instance, “*merely adjusting the weighting or coefficients is often cumbersome and not very effective*”.

2) *Motion Planner*: The description of the motion planner begins with the selection and a brief introduction to the planning algorithm. This is followed by a general description of the key components that primarily affect the performance of the planner. To gain a better understanding of how the algorithm is practically implemented, we also include the program of the key components as an additional input modality. As mentioned in Sec. I-A.2, the LLM is then able to generate repaired programs given corresponding instructions. Motivated by the chain of thought (cf. Sec. II-B), we incorporate existing explanations found within the docstrings of subfunctions to provide natural language summaries for the code blocks. The description adheres to the format of **{subfunction name}** followed by its **{docstring}**. For instance, an automatically generated **{detailed description}** is: “`self.calc_angle_to_goal` returns the orientation of the goal with respect to current position; ...” (cf. Fig. 6).

3) *Planned Trajectory*: There are various measures to quantitatively evaluate the quality of the planned trajectory and track its improvement. These measures include the cost function [56], criticality measures [62], courtesy to other traffic participants [63], and degree of traffic rule compliance [40]. To align the LLM with the desired behavior, we present not only the evaluation results for the selected measures but also incorporate the target value J^* , which can be, e.g., sourced from the motion planning benchmark leaderboard. In addition, the numerical data of the values and weights of the objective components is translated into a narrative description by mapping them to their corresponding placeholders.

4) *Few-Shots*: As it is not necessary for LLMs to have prior knowledge of the other part of the large-scale motion planner, we provide existing helper functions and their exemplary usage in the prompt. Furthermore, several human-

annotated examples for improving the performance of the specific type of motion planner can be added here, with examples available in Fig. 5.

C. LLM Querying and Iterative Prompting

When querying the LLM, it is essential to specify the desired output format. To achieve this, one can guide the LLM by emphasizing the diagnoses, prescriptions, and key components of the planner (cf. Sec. II-A) in the prompt as desired responses or employ other third-party tools such as LangChain¹. Consequently, the structured patched results can directly replace the original elements to repair the planner.

Human diagnosis experts often gauge the effectiveness of prescriptions by analyzing the outcomes of individuals and conducting follow-up consultations after the initial diagnosis. Thus, motivated by [36], [47], [64], [65], we examine the repaired planner by executing it and then pass the evaluation result back to the LLM. In case of compilation or execution errors, the previous diagnostic result is combined with the information indicating where the error occurred and what it entails, as feedback. Otherwise, the combination is made with a comparison of the performance between the updated planned trajectory and the original one.

IV. EVALUATION

We evaluate our approach using the open-source graph-search-based motion planner² from the CommonRoad platform [56], which is written in Python. As CommonRoad provides customizable challenges and annual competitions where users can compete against each other on predefined benchmarks, we can continuously integrate enhancements into DrPlanner based on insights from a broad user base. Furthermore, we choose GPT-4-Turbo³ as our LLM and use its function calling feature to generate structured outputs. The patched programs are then stringified in a JSON object and directly parsed to the motion planner, followed by execution through the `exec` function in Python. The token limit is set to 8,000, the threshold ϵ is equal to 10, and we choose the sampling temperature of the LLM at 0.6 (cf. [29, Fig. 5]). Code and exemplary prompts are available at <https://github.com/CommonRoad/drplanner>.

A. A* Search using Motion Primitives

We adapt the standard A* search algorithm using lattice-based graphs [66] (see Fig. 1), which employs a cost function and an estimated cost to the goal, namely, a *heuristic function*, to guide the search process. The graph is constructed with *motion primitives*—short trajectories generated offline through a forward simulation of a given vehicle model.

1) *Description of the Motion Planner*: The heuristic function and motion primitives constitute the key components of the planner. We provide the entire code block of the heuristic function along with descriptions of the involved subfunctions

¹<https://www.langchain.com/>

²<https://commonroad.in.tum.de/tools/commonroad-search>

³ID `gpt-4-turbo-preview` in the API of OpenAI

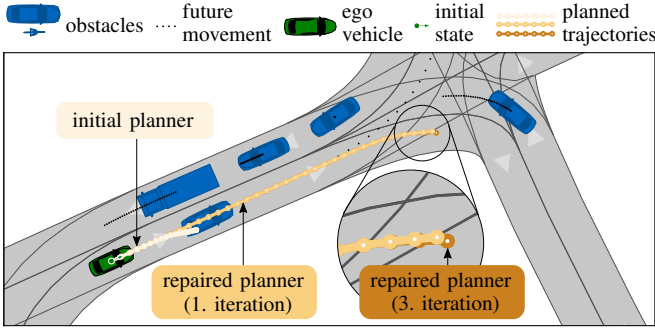


Fig. 4: Critical intersection scenario⁴ in which the ego vehicle needs to safely drive for 33 time steps. For clarity, the planned trajectories for the ego vehicle from different planners are marked with different colors and labels.

in natural language. Furthermore, motion primitives are referenced by IDs encoded with configurable parameters:

$$MP = "v_{\min} \ v_{\max} \ v_{step} \ \Delta v \ SA_{\min} \ \delta_{\max} \ S_{astep} \ \Delta \delta \ T \ \tau \ Model \ m ",$$

where v_{\min} and v_{\max} are the sampling velocity limits, δ_{\min} and δ_{\max} are the sampling steering angle bounds, Δv and $\Delta \delta$ specify their respective step sizes, τ is the time duration of each motion primitive, and m is the model identifier of the ego vehicle. All parameters are given in SI units. In the description, the explanation of the above naming convention is included, followed by the ID of motion primitives used in the original planner.

2) *Measures of the Planned Trajectory*: To evaluate the quality of the planned trajectory, we utilize the standardized objective function J_{SM1} from CommonRoad [56, Sec. VI]:

$$J_{SM1}(\chi) = \omega_A J_A + \omega_{SA} J_{SA} + \omega_{SR} J_{SR} + \omega_{LC} J_{LC} + \omega_O J_O + \omega_V J_V,$$

where ω_A , ω_{SA} , ω_{SR} , ω_{LC} , ω_O , and ω_V are the weights assigned to the respective objective components, and their values are listed in Tab. I. These components include the cost for acceleration (J_A), steering angle (J_{SA}), steering rate (J_{SR}), distance and orientation offset to the centerline of the road (J_{LC} and J_O), and velocity offset (J_V). For the evaluation scenario in Fig. 4, the target value of J_{SM1} is extracted from the CommonRoad benchmark leaderboard⁵ and is $J_{SM1}^* = 0.16$.

3) *Few-Shots*: To gain a deeper insight into the planner, we include method definitions and docstrings for existing helper functions within the planner class. As shown in Fig. 5, we also supply a list of IDs corresponding to offline-generated motion primitives, from which the LLM can select.

B. Case Study

We choose an intersection scenario from the CommonRoad platform (cf. Fig. 4), which is generated by the scenario factory for safety-critical traffic scenarios [62], [67]. In the urban environment, the search-based motion planner is

TABLE I: Comparison of the planned trajectories before and after repair. The lowest values of each objective component are marked in bold.

Item	Weight	Initial Planner	Repaired Planner	
			1. Iteration	3. Iteration
J_A	50	91.7333	14.9333	0.0000
J_{SA}	50	0.0850	0.0102	0.0147
J_{SR}	50	0.2525	0.0968	0.0673
J_{LC}	1	0.3175	0.3504	0.3393
J_O	50	0.0614	0.0038	0.0041
J_V	20	0.0000	0.0000	0.0000
J_{SM1}	-	4606.93	752.56	4.65

There are some pre-defined helper functions that can be directly called in the **heuristic function**:

```
def calc_acceleration_cost(self, path: List[KSState]) -> float:
    """Returns the acceleration costs."""
    ...
Examples:
(input)
def heuristic_function(self, node_current: PriorityNode) -> float:
    ...
    cost = angle_to_goal
    return cost
(output)
Diagnosis: the acceleration is not considered
Prescription: add the acceleration cost to the heuristic function
def heuristic_function(self, node_current: PriorityNode) -> float:
    acceleration_cost =
        self.calc_acceleration_cost(node_current.list_paths[-1])
    ...
    cost = angle_to_goal + acceleration_cost
    return cost
Feasible motion primitives with the same name format that you can
directly use:
"V_0.0_20.0_Vstep_1.0_SA_-1.066_1.066_SAstep_2.13_T_0.5_Model_BMW_320i",
"V_0.0_20.0_Vstep_2.0_SA_-1.066_1.066_SAstep_0.18_T_0.5_Model_BMW_320i",
...
```

Fig. 5: Snippet of the few-shot prompting used for the search-based planner.

responsible for navigating the ego vehicle from the initial state for 3.3s without colliding with any obstacles. The time increment of the scenario is 0.1s. Moreover, we use the planner with the setup illustrated in Fig. 6. The planned trajectory by the initial planner is shown in Fig. 4, where the ego vehicle brakes and steers slightly to the right, leading to a high value of J_{SM1} (cf. Tab. I).

The diagnostic results using our approach are illustrated in Fig. 7. In the first iteration, the provided helper functions are automatically included in the heuristic function by the LLM (cf. Fig. 7a). Meanwhile, some hyperparameters are adjusted, such as the orientation weight and the heuristic for zero velocity, and coarser motion primitives are applied. Considering all the above factors, the repaired planner results in a decrease in J_{SM1} of the planned trajectory, particularly in J_A , J_{SA} , J_{SR} , and J_O (cf. Tab. I), and leads to the vehicle traveling further forward. In contrast, the diagnostic result from the second iteration leads to a `KeyError` (cf. Fig. 7b), indicating that the repaired heuristic function is not provided by the LLM. With the iterative prompting, the error message is incorporated as feedback into the prompt for the third iteration. As shown in Fig. 7c, our approach not only helps the LLM avoid the errors from previous iterations (cf. the diagnosis “`KeyError in heuristic function`”) but also retains the previous diagnostic results that lead to a positive impact

⁴CommonRoad-ID: DEU_Guetersloh-15_2_T-1

⁵<https://commonroad.in.tum.de/solutions/ranking>

```

1 def heuristic_function(self, node_current: PriorityNode) -> float:
2   path_last = node_current.list_paths[-1]
3   angleToGoal =
4     self.calc_angle_to_goal(path_last[-1])
5   orientationToGoalDiff = self.calc_orientation_diff(angleToGoal,
6     path_last[-1].orientation)
7   cost_time = self.calc_time_cost(path_last)
8   if self.reached_goal(node_current.list_paths[-1]):
9     heur_time = 0.0
10  if self.position_desired is None:
11    heur_time = self.time_desired.start -
12      node_current.list_paths[-1][-1].time_step
13  else:
14    velocity = node_current.list_paths[-1][-1].velocity
15    if np.isclose(velocity, 0):
16      heur_time = np.inf
17    else:
18      heur_time =
19        self.calc_euclidean_distance(current_node=node_current) /
20        velocity
21  cost = 20 * orientationToGoalDiff + 0.5 * cost_time + heur_time
22  if cost < 0:
23    cost = 0
24  return cost

```

MP = "V_0.0.20.0.Vstep.4.0_SA_-1.066.1.066_SAstep.0.18.T.0.5.Model.\ BMW_320i"

Fig. 6: Heuristic function and ID of motion primitives used in the initial planner that are to be diagnosed and repaired.

on the planner. As a result, the planner significantly improves its performance, with a substantial reduction in J_{SM1} from 752.56 to 4.65, achieved by further balancing the objective components (cf. Tab. I). Moreover, it can be observed from Fig. 7 that DrPlanner can provide fine-grained diagnoses and prescriptions based on both the prompt design and fundamental aspects of programming, such as aliasing (cf. lines 10, 13, 15 in Fig. 7c). The resulting patched programs align precisely with these diagnoses and prescriptions.

C. Performance Evaluation

We further evaluate the performance of DrPlanner using the scenario illustrated in Fig. 4 and 50 A*-search-based motion planners with different setups obtained from the CommonRoad challenges. The evaluation employs the pass@ k metric and uses its unbiased version proposed in [29, Sec. 2.1], which is defined as the probability that at least one of the top $k \in \mathbb{N}_+$ generated code samples for a problem passes the given tests. Here, we use a decrease of J_{SM1} for the returned planner as the criterion for passing. Additionally, we conduct ablation studies to examine the impact of omitting two specific components within the framework: few-shots and feedback. For each study, we execute the framework 10 times to collect solution samples.

Tab. II presents the evaluation results for DrPlanner and its variants. Overall, DrPlanner exhibits robust capabilities for diagnosing and repairing motion planners, achieving a pass rate of 98.0% at $k = 10$ and an average reduction of 54.5% in J_{SM1} . Note that, similar to the case study in Sec. IV-B, the value of J_{SM1} does not converge with the iterations due to diagnostic inaccuracies. However, the average number of iterations required to observe its first decrease is 1.4. Furthermore, the results demonstrate that both the few-shot learning (cf. Sec. III-B.4) and the iterative prompting (cf. Sec. III-C) play crucial roles in enhancing the effectiveness

```

1 def heuristic_function(self, node_current: PriorityNode) -> float:
2   ...
3   acceleration_cost = self.calc_acceleration_cost(path_last)
4   path_efficiency = self.calc_path_efficiency(path_last)
5   steering_angle_cost = self.calc_steering_angle_cost(path_last)
6   steering_velocity_cost = self.calc_steering_velocity_cost(path_last)
7   if self.reached_goal(node_current.list_paths[-1]):
8     heur_time = 0.0
9     return 0
10  ...
11  if np.isclose(velocity, 0):
12    heur_time = np.inf
13  # Instead of np.inf, assign a high but finite
14  cost
15  ...
16  cost = 20 * orientationToGoalDiff + 0.5 * cost_time + heur_time
17  cost = 10 * orientationToGoalDiff + 0.5 * cost_time +
18    heur_time + acceleration_cost + path_efficiency +
19    steering_angle_cost + steering_velocity_cost
20  ...

```

MP = "V_0.0.20.0.Vstep.4.0_2.0_SA_-1.066.1.066_SAstep.0.18.T.0.5.Model.\ BMW_320i"

Diagnosis	Prescription
Orientation weight too high	Adjust the weight for orientation to goal difference
Missing acceleration cost	Include acceleration cost in heuristic
Missing path efficiency	Include path efficiency in heuristic
Missing steering angle cost	Include steering angle cost in heuristic
Missing steering velocity cost	Include steering velocity cost in heuristic
Infinite heuristic for zero velocity	Handle zero velocity case appropriately
Sparse motion primitives	Recommend motion primitives with higher branching factor

(a) First iteration.

KeyError: 'repaired_heuristic_function'

(b) Second iteration.

```

1 def heuristic_function(self, node_current: PriorityNode) -> float:
2   ...
3   acceleration_cost = self.calc_acceleration_cost(path_last)
4   path_efficiency = self.calc_path_efficiency(path_last)
5   steering_angle_cost = self.calc_steering_angle_cost(path_last)
6   steering_velocity_cost = self.calc_steering_velocity_cost(path_last)
7   if self.reached_goal(node_current.list_paths[-1]):
8     heur_time = 0.0
9     if self.position_desired is None:
10    elif self.position_desired is None:
11    heur_time = self.time_desired.start -
12      node_current.list_paths[-1][-1].time_step
13    heur_time = self.time_desired.start - path_last[-1].time_step
14  else:
15    velocity = node_current.list_paths[-1][-1].velocity
16    velocity = path_last[-1].velocity
17    if np.isclose(velocity, 0):
18      heur_time = np.inf
19    # A large but not infinite cost
20  ...
21  cost = 20 * orientationToGoalDiff + 0.5 * cost_time + heur_time
22  cost = 10 * orientationToGoalDiff + 0.5 * cost_time +
23    heur_time + acceleration_cost + path_efficiency +
24    steering_angle_cost + steering_velocity_cost
25  ...

```

MP = "V_0.0.20.0.Vstep.4.0_2.0_SA_-1.066.1.066_SAstep.0.18.T.0.5.Model.\ BMW_320i"

Diagnosis	Prescription
Orientation weight excessive	Decrease orientation weight in heuristic
Heuristic excludes costs	Incorporate all costs into heuristic
Zero velocity infinite cost	Refine zero velocity case handling
Motion primitives high branching	Optimize motion primitives branching
KeyError in heuristic function	Ensure correct key for improved heuristic

(c) Third iteration.

Fig. 7: Diagnostic and repair result for the motion planner in Fig. 6. The identical program patches in the first and third iteration are highlighted with black borders in (c). For the second iteration, we omit the diagnoses and prescriptions since it leads to an error.

TABLE II: Ablation studies on the design of DrPlanner. Values in bold denote the best performance.

Method	pass@k			Decrement of J_{SMI}	
	$k=1\uparrow$	$k=5\uparrow$	$k=10\uparrow$	Avg. \uparrow	Std. Dev. \downarrow
w/o Few-Shots	0.0%	0.0%	0.0%	0.0%	0.0%
w/o Feedback	45.4%	86.2%	92.0%	49.6%	36.3%
DrPlanner	68.0%	95.1%	98.0%	54.5%	34.9%

of DrPlanner. In particular, the few-shots prompting is more effective since the LLM is intrinsically unaware of the other supportive components of the planner, e.g., the available motion primitives. Additionally, since the initial planners are not buggy but underperforming, the results without using few-shots show that they cannot be easily improved with only the descriptions of the planner and the planned trajectory.

V. CONCLUSION

We present the first framework for diagnosing and repairing motion planners that leverages both common sense and domain-specific knowledge about causal mechanisms in LLMs. Through a modular and iterative prompt design, our approach automates the generation of descriptions for the planner and continuously enhances diagnostic performance. The major limitation of our approach is that the improvement of the planner cannot be guaranteed. However, as the capabilities of LLMs advance, we anticipate the paradigm to enhance significantly over time. Future work will involve developing datasets by monitoring user submissions over time, specifically focusing on sequences of edits that lead to performance improvements. This effort will examine the impact of different objective functions and their target values. We encourage researchers using DrPlanner to refine their motion planners and contribute towards establishing a large-scale framework that encompasses a variety of planner types for diagnostic and repair tasks.

REFERENCES

- [1] B. Paden, M. Čáp, S. Z. Yong, D. Yershov, and E. Frazzoli, “A survey of motion planning and control techniques for self-driving urban vehicles,” vol. 1, no. 1, pp. 33–55, 2016.
- [2] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, *et al.*, “Language models are few-shot learners,” *Advances in Neural Info. Processing Syst.*, vol. 33, pp. 1877–1901, 2020.
- [3] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, *et al.*, “Training language models to follow instructions with human feedback,” *Advances in Neural Info. Processing Syst.*, vol. 35, pp. 27 730–27 744, 2022.
- [4] OpenAI, “GPT-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [5] T. Gu, J. M. Dolan, and J.-W. Lee, “Runtime-bounded tunable motion planning for autonomous driving,” in *Proc. of the IEEE Intell. Veh. Symp.*, 2016, pp. 1301–1306.
- [6] S. Liu and P. Liu, “Benchmarking and optimization of robot motion planning with motion planning pipeline,” *Int. J. of Advanced Manufacturing Technology*, vol. 118, pp. 949–961, 2022.
- [7] M. Zucker, N. Ratliff, A. D. Dragan, M. Pivtoraiko, M. Klingensmith, C. M. Dellin, J. A. Bagnell, and S. S. Srinivasa, “CHOMP: Covariant Hamiltonian optimization for motion planning,” *Int. J. of Robot. Research*, vol. 32, no. 9-10, pp. 1164–1193, 2013.
- [8] S. M. LaValle, J. J. Kuffner, B. Donald, *et al.*, “Rapidly-exploring random trees: Progress and prospects,” *Algorithmic and Computational Robot.: New Directions*, vol. 5, pp. 293–308, 2001.

- [9] S. Karaman and E. Frazzoli, “Incremental sampling-based algorithms for optimal motion planning,” *Robot. Science and Systems*, vol. 104, no. 2, pp. 267–274, 2010.
- [10] S. Grigorescu, B. Trasnea, T. Cocias, and G. Macesanu, “A survey of deep learning techniques for autonomous driving,” *J. of Field Robot.*, vol. 37, pp. 362–386, 2019.
- [11] S. Aradi, “Survey of deep reinforcement learning for motion planning of autonomous vehicles,” *IEEE Trans. on Intell. Transp. Syst.*, vol. 23, no. 2, pp. 740–759, 2022.
- [12] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. A. Sallab, S. Yogamani, and P. Pérez, “Deep reinforcement learning for autonomous driving: A survey,” *IEEE Trans. on Intell. Transp. Syst.*, vol. 23, no. 6, pp. 4909–4926, 2022.
- [13] H. Krasowski, X. Wang, and M. Althoff, “Safe reinforcement learning for autonomous lane changing using set-based prediction,” in *Proc. of the IEEE Int. Conf. on Intell. Transp. Syst.*, 2020, pp. 1–7.
- [14] L. Wang, L. Sun, M. Tomizuka, and W. Zhan, “Socially-compatible behavior design of autonomous vehicles with verification on real human data,” *IEEE Robot. and Automation Letters*, vol. 6, no. 2, pp. 3421–3428, 2021.
- [15] N. Mehdipour, M. Althoff, R. D. Tebbens, and C. Belta, “Formal methods to comply with rules of the road in autonomous driving: State of the art and grand challenges,” *Automatica*, vol. 152, no. 110692, 2023.
- [16] L. Gazzola, D. Micucci, and L. Mariani, “Automatic software repair: A survey,” *IEEE Trans. on Software Engineering*, vol. 45, no. 1, pp. 34–67, 2019.
- [17] W. Weimer, T. Nguyen, C. Le Goues, and S. Forrest, “Automatically finding patches using genetic programming,” in *Proc. of the IEEE Int. Conf. on Software Engineering*, 2009, pp. 364–374.
- [18] D. Kim, J. Nam, J. Song, and S. Kim, “Automatic patch generation learned from human-written patches,” in *Proc. of the IEEE Int. Conf. on Software Engineering*, 2013, pp. 802–811.
- [19] X. Liu and H. Zhong, “Mining StackOverflow for program repair,” in *Proc. of the IEEE Int. Conf. on Software analysis, Evolution and Reengineering*, 2018, pp. 118–129.
- [20] K. Liu, A. Koyuncu, D. Kim, and T. F. Bisseyandé, “TBar: Revisiting template-based automated program repair,” in *Proc. of the ACM Int. Symp. on Software Testing and Analysis*, 2019, pp. 31–42.
- [21] A. Koyuncu, K. Liu, T. F. Bisseyandé, D. Kim, J. Klein, M. Monperrus, and Y. Le Traon, “Fixminer: Mining relevant fix patterns for automated program repair,” *Empirical Software Engineering*, vol. 25, pp. 1980–2024, 2020.
- [22] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” in *Int. Conf. on Learning Representations*, 2015.
- [23] T. Lutellier, H. V. Pham, L. Pang, Y. Li, M. Wei, and L. Tan, “CoCoNut: combining context-aware neural translation models using ensemble for program repair,” in *Proc. of the ACM Int. Symp. on Software Testing and Analysis*, 2020, pp. 101–114.
- [24] Q. Zhu, Z. Sun, Y.-a. Xiao, W. Zhang, K. Yuan, Y. Xiong, and L. Zhang, “A syntax-guided edit decoder for neural program repair,” in *Proc. of the ACM Joint Meeting on European Software Engineering Conf. and Symp. on the Foundations of Software Engineering*, 2021, pp. 341–353.
- [25] N. Jiang, T. Lutellier, and L. Tan, “CURE: Code-aware neural machine translation for automatic program repair,” in *Proc. of the IEEE/ACM Int. Conf. on Software Engineering*, 2021, pp. 1161–1173.
- [26] H. Ye, M. Martinez, and M. Monperrus, “Neural program repair with execution-based backpropagation,” in *Proc. of the IEEE/ACM Int. Conf. on Software Engineering*, 2022, pp. 1506–1518.
- [27] C. S. Xia and L. Zhang, “Less training, more repairing please: revisiting automated program repair via zero-shot learning,” in *Proc. of the ACM Joint European Software Engineering Conf. and Symp. on the Foundations of Software Engineering*, 2022, pp. 959–971.
- [28] Z. Feng, D. Guo, D. Tang, N. Duan, X. Feng, M. Gong, L. Shou, B. Qin, T. Liu, D. Jiang, *et al.*, “CodeBERT: A pre-trained model for programming and natural languages,” in *Proc. of the Conf. on Empirical Methods in Natural Language Processing*, 2020, pp. 1536–1547.
- [29] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. d. O. Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, *et al.*, “Evaluating large language models trained on code,” *arXiv preprint arXiv:2107.03374*, 2021.

- [30] J. Austin, A. Odena, M. Nye, M. Bosma, H. Michalewski, D. Dohan, E. Jiang, C. Cai, M. Terry, Q. Le, *et al.*, “Program synthesis with large language models,” *arXiv preprint arXiv:2108.07732*, 2021.
- [31] D. Fried, A. Aghajanyan, J. Lin, S. Wang, E. Wallace, F. Shi, R. Zhong, W.-t. Yih, L. Zettlemoyer, and M. Lewis, “InCoder: A generative model for code infilling and synthesis,” in *Proc. of the Int. Conf. on Learning Representations*, 2023.
- [32] S. D. Kolak, R. Martins, C. Le Goues, and V. J. Hellendoorn, “Patch generation with language models: Feasibility and scaling behavior,” in *Proc. of the Int. Conf. on Learning Representations: Deep Learning for Code Workshop*, 2022.
- [33] J. A. Prenner, H. Babii, and R. Robbes, “Can OpenAI’s Codex fix bugs? an evaluation on QuixBugs,” in *Proc. of the Int. Workshop on Automated Program Repair*, 2022, pp. 69–75.
- [34] C. S. Xia, Y. Wei, and L. Zhang, “Practical program repair in the era of large pre-trained language models,” in *Proc. of the Int. Conf. on Software Engineering*, 2023, pp. 1482–1494.
- [35] N. Shinn, F. Cassano, A. Gopinath, K. R. Narasimhan, and S. Yao, “Reflexion: Language agents with verbal reinforcement learning,” *Advances in Neural Info. Processing Syst.*, vol. 36, 2024.
- [36] X. Chen, M. Lin, N. Schärli, and D. Zhou, “Teaching large language models to self-debug,” in *Proc. of the Int. Conf. on Learning Representations*, 2023.
- [37] N. Jain, T. Zhang, W.-L. Chiang, J. E. Gonzalez, K. Sen, and I. Stoica, “LLM-assisted code cleaning for training accurate code generators,” *arXiv preprint arXiv:2311.14904*, 2023.
- [38] A. Madaan, A. Shypula, U. Alon, M. Hashemi, P. Ranganathan, Y. Yang, G. Neubig, and A. Yazdanbakhsh, “Learning performance-improving code edits,” in *Proc. of the Int. Conf. on Learning Representations*, 2024.
- [39] Y. Lin, S. Maierhofer, and M. Althoff, “Sampling-based trajectory repairing for autonomous vehicles,” in *Proc. of the IEEE Int. Conf. on Intell. Transp. Syst.*, 2021, pp. 572–579.
- [40] Y. Lin and M. Althoff, “Rule-compliant trajectory repairing using satisfiability modulo theories,” in *Proc. of the IEEE Intell. Veh. Symp.*, 2022, pp. 449–456.
- [41] S. Maierhofer, Y. Ballnath, and M. Althoff, “Map verification and repairing using formalized map specifications,” in *Proc. of the IEEE Int. Conf. on Int. Transp. Syst.*, 2023, pp. 1277–1284.
- [42] A. Pacheck and H. Kress-Gazit, “Physically feasible repair of reactive, linear temporal logic-based, high-level tasks,” *IEEE Trans. on Robot.*, vol. 39, no. 6, pp. 4653–4670, 2023.
- [43] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, “Large language models are zero-shot reasoners,” *Advances in Neural Info. Processing Syst.*, vol. 35, pp. 22 199–22 213, 2022.
- [44] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafraan, K. Narasimhan, and Y. Cao, “ReAct: Synergizing reasoning and acting in language models,” in *Proc. of the Int. Conf. on Learning Representations*, 2023.
- [45] E. Kiciman, R. Ness, A. Sharma, and C. Tan, “Causal reasoning and large language models: Opening a new frontier for causality,” *arXiv preprint arXiv:2305.00050*, 2023.
- [46] H. Sha, Y. Mu, Y. Jiang, L. Chen, C. Xu, P. Luo, S. E. Li, M. Tomizuka, W. Zhan, and M. Ding, “LanguageMPC: Large language models as decision makers for autonomous driving,” *arXiv preprint arXiv:2310.03026*, 2023.
- [47] L. Wen, D. Fu, X. Li, X. Cai, T. Ma, P. Cai, M. Dou, B. Shi, L. He, and Y. Qiao, “DiLu: A knowledge-driven approach to autonomous driving with large language models,” in *Proc. of the Int. Conf. on Learning Representations*, 2024.
- [48] W. Wang, J. Xie, C. Hu, H. Zou, J. Fan, W. Tong, Y. Wen, S. Wu, H. Deng, Z. Li, *et al.*, “DriveMLM: Aligning multi-modal large language models with behavioral planning states for autonomous driving,” *arXiv preprint arXiv:2312.09245*, 2023.
- [49] C. Sima, K. Renz, K. Chitta, L. Chen, H. Zhang, C. Xie, P. Luo, A. Geiger, and H. Li, “DriveLM: Driving with graph visual question answering,” *arXiv preprint arXiv:2312.14150*, 2023.
- [50] C. Cui, Y. Ma, X. Cao, W. Ye, and Z. Wang, “Drive as you speak: Enabling human-like interaction with large language models in autonomous vehicles,” in *Proc. of the IEEE/CVF Winter Conf. on Applications of Computer Vision*, 2024, pp. 902–909.
- [51] J. Mao, Y. Qian, H. Zhao, and Y. Wang, “GPT-driver: Learning to drive with GPT,” *arXiv preprint arXiv:2310.01415*, 2023.
- [52] J. Mao, J. Ye, Y. Qian, M. Pavone, and Y. Wang, “A language agent for autonomous driving,” *arXiv preprint arXiv:2311.10813*, 2023.
- [53] Z. Xu, Y. Zhang, E. Xie, Z. Zhao, Y. Guo, K. K. Wong, Z. Li, and H. Zhao, “DriveGPT4: Interpretable end-to-end autonomous driving via large language model,” *arXiv preprint arXiv:2310.01412*, 2023.
- [54] L. Chen, O. Sinavski, J. Hünemann, A. Karnsund, A. J. Willmott, D. Birch, D. Maund, and J. Shotton, “Driving with LLMs: Fusing object-level vector modality for explainable autonomous driving,” *arXiv preprint arXiv:2310.01957*, 2023.
- [55] H. Shao, Y. Hu, L. Wang, S. L. Waslander, Y. Liu, and H. Li, “LM-Drive: Closed-loop end-to-end driving with large language models,” *arXiv preprint arXiv:2312.07488*, 2023.
- [56] M. Althoff, M. Koschi, and S. Manzingler, “CommonRoad: Composable benchmarks for motion planning on roads,” in *Proc. of the IEEE Intell. Veh. Symp.*, 2017, pp. 719–726.
- [57] C. Pek, V. Rusinov, S. Manzingler, M. C. Üste, and M. Althoff, “CommonRoad Drivability Checker: Simplifying the development and validation of motion planning algorithms,” in *Proc. of the IEEE Intell. Veh. Symp.*, 2020, pp. 1013–1020.
- [58] P. Liu, W. Yuan, J. Fu, Z. Jiang, H. Hayashi, and G. Neubig, “Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing,” *ACM Computing Surveys*, vol. 55, no. 9, pp. 1–35, 2023.
- [59] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, *et al.*, “Chain-of-thought prompting elicits reasoning in large language models,” *Advances in Neural Info. Processing Syst.*, vol. 35, pp. 24 824–24 837, 2022.
- [60] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in Neural Info. Processing Syst.*, vol. 30, pp. 5998–6008, 2017.
- [61] J. Yang, H. Jin, R. Tang, X. Han, Q. Feng, H. Jiang, B. Yin, and X. Hu, “Harnessing the power of LLMs in practice: A survey on ChatGPT and beyond,” *arXiv preprint arXiv:2304.13712*, 2023.
- [62] Y. Lin and M. Althoff, “CommonRoad-CriMe: A toolbox for criticality measures of autonomous vehicles,” in *Proc. of the IEEE Intell. Veh. Symp.*, 2023, pp. 1–8.
- [63] W. Schwarting, A. Pierson, J. Alonso-Mora, S. Karaman, and D. Rus, “Social behavior for autonomous vehicles,” *Proc. of the National Academy of Sciences*, vol. 116, no. 50, pp. 24 972–24 978, 2019.
- [64] Z. Liu, A. Bahety, and S. Song, “REFLECT: Summarizing robot experiences for failure explanation and correction,” in *Proc. of the Conference on Robot Learning*, 2023.
- [65] M. Skreta, N. Yoshikawa, S. Arellano-Rubach, Z. Ji, L. B. Kristensen, K. Darvish, A. Aspuru-Guzik, F. Shkurti, and A. Garg, “Errors are useful prompts: Instruction guided task programming with verifier-assisted iterative prompting,” *arXiv preprint arXiv:2303.14100*, 2023.
- [66] M. Pivtoraiko and A. Kelly, “Kinodynamic motion planning with state lattice motion primitives,” in *Proc. of the IEEE/RSJ Int. Conf. on Intell. Robots and Syst.*, 2011, pp. 2172–2179.
- [67] M. Klischat, E. I. Liu, F. Holtke, and M. Althoff, “Scenario factory: Creating safety-critical traffic scenarios for automated vehicles,” in *Proc. of the IEEE Int. Conf. on Intell. Transp. Syst.*, 2020, pp. 1–7.