# Trajectory-guided Control Prediction for End-to-end Autonomous Driving: A Simple yet Strong Baseline

**Penghao Wu**[*][†]
Shanghai AI Laboratory
Shanghai Jiao Tong University
wupenghaocraig@sjtu.edu.cn

**Xiaosong Jia**[*]
Shanghai Jiao Tong University
Shanghai AI Laboratory
jiaxiaosong@sjtu.edu.cn

**Li Chen**[*]
Shanghai AI Laboratory
lichen@pjlab.org.cn

**Junchi Yan**[†]
Shanghai Jiao Tong University
Shanghai AI Laboratory
yanjunchi@sjtu.edu.cn

**Hongyang Li**
Shanghai AI Laboratory
Shanghai Jiao Tong University
lihongyang@pjlab.org.cn

**Yu Qiao**
Shanghai AI Laboratory
qiaoyu@pjlab.org.cn

## Abstract

Current end-to-end autonomous driving methods either run a controller based on a planned trajectory or perform control prediction directly, which have spanned two separately studied lines of research. Seeing their potential mutual benefits to each other, this paper takes the initiative to explore the combination of these two well-developed worlds. Specifically, our integrated approach has two branches for trajectory planning and direct control, respectively. The trajectory branch predicts the future trajectory, while the control branch involves a novel multi-step prediction scheme such that the relationship between current actions and future states can be reasoned. The two branches are connected so that the control branch receives corresponding guidance from the trajectory branch at each time step. The outputs from two branches are then fused to achieve complementary advantages. Our results are evaluated in the closed-loop urban driving setting with challenging scenarios using the CARLA simulator. Even with a monocular camera input, the proposed approach ranks *first* on the official CARLA Leaderboard, outperforming other complex candidates with multiple sensors or fusion mechanisms by a large margin. The source code is publicly available at https://github.com/OpenPerceptionX/TCP.

## 1 Introduction

End-to-end autonomous driving methods, which directly map raw sensor data to a planned trajectory or low-level control actions, show the virtue of simplicity, conceptually avoiding the cascading error of complex modular design and heavy hand-crafted rules. The output prediction of the model for end-to-end autonomous driving generally falls into two forms: trajectory/waypoints [48, 4, 11, 46, 15, 29, 10] and direct control actions [17, 39, 18, 42, 12, 60, 9]. However, there is still no clear conclusion as to which of these two forms is better for all circumstances or certain scenarios.

Different from control predictions that could be directly applied to the vehicle, for methods that **plan trajectory**, additional controllers such as PID controllers are usually needed as a subsequent step to convert the planned trajectory into control signals. One attractive and potential supremacy of trajectory-based prediction is that it actually considers a relatively longer time horizon into the future and could be further combined with other modules (*e.g.*, multi-agent trajectory prediction [59, 10],

---

[*] Equal Contribution. Work done when PW and XJ were interns at Shanghai AI Laboratory.
[†] Correspondence author.

Figure 1: Typical failure cases of two prediction paradigms. Red dots indicate the trajectory prediction, blue dots are the actual path following the red trajectory with PID controllers, and green dots denote the actual path from control-based method. (a) trajectory-based methods may struggle for big turns. (b) control-based methods may have a reaction latency and suffer from abrupt obstacles due to focusing on current time step only. These observations motivate us to propose a unified framework to combine these two worlds for mutual benefits.

semantic or occupancy prediction modules [15, 8, 25]) to reduce possible collisions. However, turning the trajectory into control actions so that the vehicle could follow the planned trajectory is not trivial [57]. The industry usually adopts sophisticated control algorithms such as model predictive control to achieve reliable trajectory-following performance [5, 21]. Simple PID controllers may perform worse in situations such as taking a big turn or starting at the red light due to the inertial problem of end-to-end models [29]. For **control-based** methods, the control signals are directly optimized. Nevertheless, their focus on the current step may cause deferred reactions to avoid potential collisions with other moving agents. The independence between the control predictions of different steps also makes the actions of the vehicle more unstable or discontinuous. Fig. 1 shows two typical cases where two paradigms fail respectively. How to combine these two forms of prediction model as well as their outputs is an interesting yet relatively rarely studied area, which motivates this work.

One straightforward (but in fact rarely studied in literature) idea is to train a control prediction model and a trajectory planning model separately, and combine their ultimate outputs directly. It can be viewed as an ensemble of two different models. However, such a naive approach not only doubles the size of the model, but also ignores possible useful correlations between these two forms. To this end, we introduce the **TCP** (Trajectory-guided Control Prediction) framework, packing these two branches into a unified framework. It can be viewed as a multi-task learning (MTL) [7, 2] framework where a shared backbone extracts common features with decreased computational complexity as well as the increased ability of generalization due to the close relationship between the two tasks [38, 34, 13]. Furthermore, to address the drawbacks of current control prediction methods, we delicately devise a novel multi-step control branch and a trajectory-guided control prediction scheme.

While trajectory planning considers several steps into the future, directly learning the control in a behavior cloning fashion [44, 41, 17, 18, 11] often focuses on the current time step only, given prior on each state-action pair as independent and identical distributed (IID). This assumption is not accurate and may hamper the long-term performance since the driving task is a sequential decision-making problem. To alleviate the problem, we propose to predict multi-step control actions into the future. However, the multi-step control process needs interactions with the environment. Thus we formulate a temporal module to learn the forward process and interactions between the ego agent and the environment. A temporal module implemented with GRU [16] progressively deals with the feature representation for each time step, implicitly taking into account the dynamic motion of agents, interaction among them and dynamic environment information such as the changing of traffic lights.

Additionally, to generate accurate control signals in the multi-step prediction scheme, the model should retrieve proper location information from current sensor input for different future time steps. For example, an agent may pay more attention to nearby regions for a few early future time steps and far away regions for the remote ones. Considering that the knowledge has already been partly encoded in the trajectory branch, we adopt the attention mechanism to locate those critical and helpful areas in the long-term trajectory prediction branch, and guide the control prediction branch to pay attention to them at each future step in a corresponding way. As a result, our model is capable of reasoning about how to optimize current control prediction so that the future states are similar to those from the expert when the predicted control actions are applied.

2

With the predicted trajectory and control signals from two branches, we propose a situation based fusion scheme to adaptively combine these two forms in a self-ensemble way to form the ultimate output according to the experiments results and prior knowledge. It combines the best of these two forms, which further boosts the performance under different scenarios.

TCP has shown superior performance when being validated in the CARLA driving simulator [20]. Our method, which only uses a **monocular camera**, achieves a **75.137** driving score and ranks **1st** on the public CARLA Leaderboard [1], even surpassing prior state-of-the-art methods using multiple cameras and a LiDAR by 13.291 points. The main **contributions** of this paper include:

- We examine two dominant paradigms for end-to-end autonomous driving: trajectory planning and direct control, and propose to combine them in an integrated learning pipeline. To our knowledge, this is the first time that such two branches are jointly learned and fused for prediction.

- A multi-step control prediction branch with a temporal module and trajectory-guided attention is devised to enable temporal reasoning. To combine the best of two branches, we design a situation based scheme to fuse the two outputs.

- As a simple yet strong baseline, our method with only a monocular camera as input achieves new state-of-the-art on the CARLA Leaderboard with many competitors using multiple sensors. We conduct thorough ablation studies to verify the effectiveness of our approach.

## 2 Related Work

### 2.1 End-to-end Autonomous Driving

Learning-based end-to-end autonomous driving has emerged as an active research topic in recent years. Studies usually fall into two categories: reinforcement learning (RL) and imitation learning. RL is a promising way to address the problem of being more robust to the distribution shifts of datasets. Liang *et al.* [39] use DDPG to train a policy which is pre-trained in a supervised way. Kendall *et al.* [30] train their deep RL algorithm onboard to efficiently learn to drive a real-world vehicle. The perception task is decoupled out of the online RL process in [50, 9, 62]. The model-based method WoR [12] assumes world on rails and uses policy distillation to realize powerful performance.

Imitation learning (IL), especially behavior cloning, collects recorded data for models to mimic with high data efficiency. The expert data typically has two forms, trajectories and control actions. Zeng *et al.* [58] train a cost volume to generate the planning route, while [49, 8, 25] explicitly design safety and comfort costs based on semantic occupancy maps to select the best one in the expert trajectory sets. Zhang *et al.* [59] predict trajectories of surrounding vehicles with labeled BEV map. LBC [11] and NEAT [15] decode waypoints from a dense heatmap or offset map. These approaches aforementioned all utilize a relatively dense representation to obtain results which increases model complexity. Transfuser and its variants [46, 29] adopt a simple GRU to auto-regress waypoints. LAV [10] adopts a temporal GRU module to further refine the trajectory. They unanimously achieve impressive performance on the CARLA leaderboard, motivating us to adapt the auto-regression scheme as well in our design. On the other hand, all trajectory-based methods use PID controllers to get the ultimate actions, which may cause inferior effects in complicated scenarios.

Another genre to predict control actions directly is proposed in [44, 40, 3, 54]. CIL [17] adds a measurement encoder and multiple branches for different high level commands with the image encoder. CILRS [18] is proposed afterwards and further introduces a speed prediction head. They stand as classic baselines for IL in the CARLA driving simulator. Diverse optimized approaches are presented based on them, such as multi-modal inputs [22, 53], multi-task learning [56, 37, 24, 27, 31, 26, 63], dataset aggregation [45] and knowledge distillation [61, 60]. However, the compact control-based methods often have higher vehicles collision rates, remaining an interesting domain to explore. Similar work exists in other related domains such as robotic navigation as well. [43] learns a controller after a local trajectory planner to improve the overall navigation behavior.

### 2.2 Multi-task and Ensemble Learning for Autonomous Driving

Multi-task learning is a popular approach to train several related tasks simultaneously to help each other and improve generalization [7, 2]. Combinations of various autonomous driving tasks such

as object detection, lane detection, semantic segmentation, depth estimation, *etc.* have been proved to be capable of achieving incredible performance [38, 14, 47, 34, 52, 13]. MTL is also suitable in the end-to-end problem since it is observed the performance of a direct mapping from an image to control signals is limited. [56] adds a speed prediction task similar to CIL [17] and [63] separates the lateral and longitudinal controls as two tasks. LAV [10] trains an extra scene mapping network, and [24, 27, 29] additionally predict optical flow or dense depth. Our idea of training trajectory and control simultaneously is closely related to FASNet [31]. FASNet predicts future positions of the ego agent as an auxiliary task and adds a kinematic loss considering the relation between control and locations. However, the constrain is based on a constant velocity model which neglects the important throttle and brake, and it does not work at the inference time. On the other hand, our TCP framework has feature interactions at an earlier stage to fully explore their potential mutual benefits.

Ensembles of models have long been utilized to improve the performance in computer vision [19, 33, 35, 51, 55, 45]. Besides the normal combination of models, two classic ensemble learning methods are particularly preferred in the autonomous driving regime. One is the Test-Time Augmentation (TTA), which is of great help to the 3D object detection task with LiDAR [6, 36]. Another one is the fusion of experts [28] where experts are trained on a subset of the input space and a gating network is trained to provide the fusion weights. LSD [42] and MoDE [32] divide a dataset into sub-scenarios to get different sub-policies for end-to-end autonomous driving. These traditional ensemble approaches combine models of the same structure while our approach tries to combine two different representations. Also, the multiple experts design increases the complexity of the training strategy and we seek to have a simpler situation based fusion scheme to boost the performance.

## 3 Trajectory-guided Control Prediction

### 3.1 Problem Setting

**Problem formulation.** Given the state $\mathbf{x}$ comprised of the sensor signal $\mathbf{i}$, the speed of the vehicle $v$, and the high level navigation information $\mathbf{g}$ including a discrete navigation command and the coordinates of navigation target provided by the global planner, the end-to-end model needs to output control signals $\mathbf{a}$ comprised of longitudinal control signals *throttle* $\in [0, 1]$ and *brake* $\in [0, 1]$, and the lateral control signal *steer* $\in [-1, 1]$.

Conventional methods tackle this problem with either a trajectory-output or a control-output only model. However, TCP combines both of them as two branches: a trajectory branch which predicts the planned trajectory and a control branch which is guided by the trajectory one and outputs both current and multi-step control signals into the future. Both branches are trained in a supervised manner. Consider an expert which directly outputs the control signals at each step, supervising the predicted trajectory with the ground truth trajectory makes it not strictly satisfy the setting of behavior cloning in imitation learning. The ground truth trajectory indeed involves future expert actions and future states about the environment, so we formulate it as a trajectory planning task with ground truth trajectory as supervision for our trajectory branch. As for the control branch, training a control model which makes current control prediction supervised by the expert control is just behavior cloning in imitation learning, and it can be formulated as:

$$\arg \min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{a}^*) \sim D}[\mathcal{L}(\mathbf{a}^*, \pi_\theta(\mathbf{x}))], \tag{1}$$

where $D = \{(\mathbf{x}, \mathbf{a}^*)\}$ is a dataset comprised of state-action pairs collected from the expert. $\pi_\theta$ denotes the policy of the control branch, and $\mathcal{L}$ is the loss measuring how close the action from the expert and the action from our model is. The expert collects the dataset by controlling the vehicle and interacting with the world. Each collected route is a trajectory $\xi = (\mathbf{x}_0, \mathbf{a}_0^*, \mathbf{x}_1, \mathbf{a}_1^*, \cdots, \mathbf{x}_T)$ as a sequence of state action pairs $\{(\mathbf{x}_i, \mathbf{a}_i^*)\}_{i=0}^{T}$, which is then added into the whole dataset $D$.

**Expert demonstration.** Here we choose Roach [60] as the expert. Roach is a simple model trained by RL with privileged information, including roads, lanes, routes, vehicles, pedestrians, traffic lights, and stops, all being rendered into a 2D BEV image. Such a learning-based expert can transfer more information besides the direct supervision signals compared with an expert made by hand-crafted rules. Specifically, we have a feature loss which forces the latent features before the final output head from the student model to be similar to that of the expert. A value loss is also added as an auxiliary task for the student model to predict an expected return.
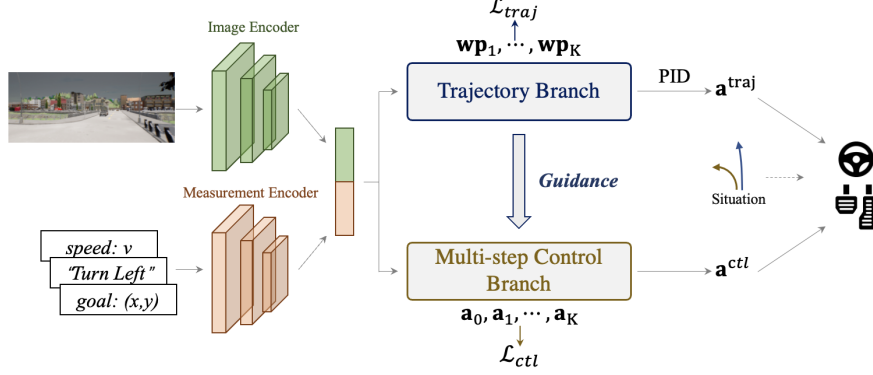
Figure 2: Overview of Trajectory-guided Control Prediction (TCP). The encoded features are shared by the trajectory and multi-step control branch. The trajectory branch provides per-step guidance for multi-step control prediction. Outputs from two branches are combined according to our situation based fusion scheme to generate the ultimate control actions.

## 3.2 Architecture Design

**Overview.** As illustrated in Fig. 2, the whole architecture is comprised of an input encoding stage and two subsequent branches. The input image $\mathbf{i}$ goes through a CNN based image encoder, such as ResNet [23], to generate a feature map $\mathbf{F}$. In the meantime, the navigation information $\mathbf{g}$ is concatenated with the current speed $v$ to form the measurement input $\mathbf{m}$, then an MLP based measurement encoder takes $\mathbf{m}$ as its input and outputs the measurement feature $\mathbf{j}_{\mathrm{m}}$. The encoded features are then shared by two branches for subsequent trajectory and control predictions. Specifically, the control branch is a novel multi-step prediction design with guidance from the trajectory one, which will be illustrated in detail in the following sections. Finally, a situation based fusion scheme is adopted to combine the best of the two output paradigms. We will go over each part in detail below.

### 3.2.1 Trajectory planning branch

Different from control prediction which directly predicts control actions, the trajectory planning branch first generates a planned trajectory comprised of waypoints at $K$ steps for the agent to follow, and then the trajectory is processed by low-level controllers to get the final control actions. With the shared feature from the input encoder, the image feature map $\mathbf{F}$ is average pooled and concatenated with the measurement feature $\mathbf{j}_{\mathrm{m}}$ to form $\mathbf{j}^{\mathrm{traj}}$. Inspired by [46], we feed $\mathbf{j}^{\mathrm{traj}}$ into a GRU [16] to auto-regressively obtain future waypoints one by one to form the planned trajectory altogether.

We have two PID controllers for longitudinal and lateral control respectively. With the planned trajectory, we first calculate the vectors between consecutive waypoints. The magnitudes of these vectors represent the desired speed and are sent to the longitudinal controller to generate $throttle$ and $brake$ control actions, and the orientations are sent to the lateral controller to get the $steer$ action.

### 3.2.2 Multi-step control prediction branch

As discussed in Sec. 3.1, for a control model predicting current control actions based on current input only, the supervised training is just behavior cloning, which relies on the independent and identically distributed (IID) assumption. This assumption apparently does not hold because of the distribution shifts in test cases, since the closed-loop tests require sequential decision making where the historical actions will affect the future states and actions. Instead of modeling it as a Markov Decision Process (MDP) and resorting to reinforcement learning, here we devise a simple way to mitigate the problem by predicting multi-step control into the future.

Given the current state $\mathbf{x}_t$, now our multi-step control prediction branch outputs multiple actions: $\pi_{\theta_{multi}} = (\mathbf{a}_t, \mathbf{a}_{t+1}, \cdots, \mathbf{a}_{t+K})$. However, it is difficult to predict future control actions since we only have sensor inputs at the current time step. Towards this problem, we devise a temporal module to implicitly carry out the changing and interaction process of the environment and our agent. It is supposed to provide mainly dynamic information about the environment and the status of the agent

itself, such as the motion of other objects, the changing of traffic lights, and the status of the ego agent. Meanwhile, to improve the ability of incorporating critical static information (*e.g.*, curbs and lanes) and boost the spatial consistency of two branches, we propose to use the trajectory branch to guide the control counterpart to attend to proper regions of the input image at each future time step.

**Temporal module.** Our temporal module is implemented with a GRU for better consistency with the trajectory branch. At step $t$ ($0 \leq t \leq K - 1$), the input for the temporal module is the concatenation of the current feature $\mathbf{j}_t^{\mathrm{ctl}}$ (more construction details in the next section) and current predicted action $\mathbf{a}_t$, which is a compact representation about the current states of the environment and the agent itself. The temporal module is supposed to reason about the dynamic changing process based on current feature vector and the predicted action. Then the updated hidden state $\mathbf{h}_{t+1}^{\mathrm{ctl}}$ will contain dynamic information about the environment and the updated status of the agent at time step $t + 1$. To some extent, the temporal module acts as a coarse simulator with the whole environment and the agent being abstracted as a feature vector. It then simulates the interaction between the environment and the agent based on current prediction of actions.

**Trajectory-guided attention.** With the sensor input at current step only, it is hard to pick out desirable regions where the model should focus on at future steps. However, the location of the ego agent contains important cues about how to find those regions containing critical static information for control prediction at each step.

Therefore, we seek help from the trajectory planning branch to get information about the possible location of our agent at that corresponding step. As shown in Fig. 3, TCP implements this by learning an attention map to extract important information from the encoded feature map. The interaction between two branches enhances the consistency of these two strongly related output paradigms and further elaborates the multi-task spirit. Specifically, with the 2D feature map extracted by the image encoder $\mathbf{F}$ at time step $t$ ($1 \leq t \leq K$), we calculate an attention map $\mathbf{w}_t \in \mathbb{R}^{1 \times \mathrm{H} \times \mathrm{W}}$ using the corresponding hidden states from the control branch and the trajectory branch:



Figure 3: Detailed trajectory guiding process. For predictions at time step $t$, the hidden states from the waypoint GRU and the temporal module are combined to learn an attention weight map to re-aggregate the 2D image feature map for control prediction.

$$\mathbf{w}_t = \mathrm{MLP}(\mathrm{Concat}[\mathbf{h}_t^{\mathrm{traj}}, \mathbf{h}_t^{\mathrm{ctl}}]). \tag{2}$$

The attention map $\mathbf{w}_t \in \mathbb{R}^{1 \times \mathrm{H} \times \mathrm{W}}$ is adopted to aggregate the feature map $\mathbf{F}$ for this step. We then combine the attended feature map with $\mathbf{h}_t^{\mathrm{ctl}}$ to form the informative representation feature $\mathbf{j}_t^{\mathrm{ctl}}$ containing both static and dynamic information about the environment and the ego agent. The process can be described as follows:

$$\mathbf{j}_t^{\mathrm{ctl}} = \mathrm{MLP}(\mathrm{Concat}[\mathrm{Sum}(\mathrm{Softmax}(\mathbf{w}_t) \odot \mathbf{F}), \mathbf{h}_t^{\mathrm{ctl}}]). \tag{3}$$

The informative representation feature $\mathbf{j}_t^{\mathrm{ctl}}$ is fed into a policy head which is shared among all time steps to predict the corresponding control action $\mathbf{a}_t$. Note that for the initial step, we only use the measurement feature to calculate the initial attention map and combine the attended image feature with the measurement feature to form the initial feature vector $\mathbf{j}_0^{\mathrm{ctl}}$. To guarantee the feature $\mathbf{j}_t^{\mathrm{ctl}}$ does describe the state at that step and contain the important information for control prediction, we add a feature loss at each step to make $\mathbf{j}_t^{\mathrm{ctl}}$ close to the feature of the expert as well.

To this end, our TCP framework endows the model with the reasoning ability along a short time horizon. It emphasizes how to make current control prediction close to the one from the expert. Furthermore, it takes into account what current control prediction can make the environment states and status of ego agent in future time steps similar to the ones from the expert.

### 3.3 Loss Design

Our loss contains trajectory planning loss $\mathcal{L}_{traj}$, control prediction $\mathcal{L}_{ctl}$, and auxiliary loss $\mathcal{L}_{aux}$.
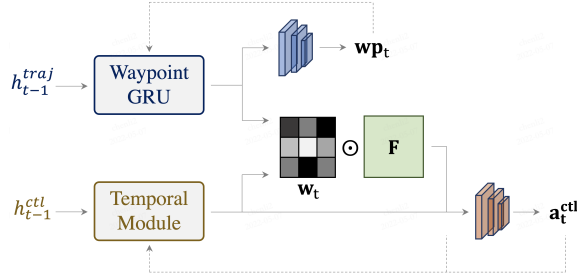
For the trajectory planning branch, the loss $\mathcal{L}_{traj}$ can be expressed as:

$$\mathcal{L}_{traj} = \sum_{t=1}^{K} \|\mathbf{wp}_t - \mathbf{w\hat{p}}_t\|_1 + \lambda_F \cdot \mathcal{L}_F \left( \mathbf{j}_0^{\text{traj}}, \mathbf{j}_0^{\text{Expert}} \right), \tag{4}$$

where $\mathbf{wp}_t$, $\mathbf{w\hat{p}}_t$ are the predicted and ground truth waypoint at the $t^{th}$ step respectively. $\mathcal{L}_F$ indicates the feature loss measuring the $L_2$ distance between $\mathbf{j}_0^{\text{traj}}$ and the feature $\mathbf{j}_0^{\text{Expert}}$ from the expert at the current step as an additional supervision signal [60]. $\lambda_F$ is a tunable loss weight.

For the control prediction branch, we model the action as a beta distribution. The loss $\mathcal{L}_{ctl}$ is:

$$\begin{aligned}
\mathcal{L}_{ctl} =& \mathbf{KL} \left( \text{Beta}(\mathbf{a}_0) || \text{Beta}(\mathbf{\hat{a}}_0) \right) + \frac{1}{K} \sum_{t=1}^{K} \mathbf{KL} \left( \text{Beta}(\mathbf{a}_t) || \text{Beta}(\mathbf{\hat{a}}_t) \right) \\
& + \lambda_F \cdot \mathcal{L}_F \left( \mathbf{j}_0^{\text{ctl}}, \mathbf{j}_0^{\text{Expert}} \right) + \frac{1}{K} \sum_{t=1}^{K} \mathcal{L}_F \left( \mathbf{j}_t^{\text{ctl}}, \mathbf{j}_t^{\text{Expert}} \right),
\end{aligned} \tag{5}$$

where $\text{Beta}(\mathbf{a})$ denotes the beta distribution represented by the corresponding predicted distribution parameters and KL-divergence is used to measure the similarity between the predicted control distribution and the one from expert, *i.e.*, $\text{Beta}(\mathbf{\hat{a}})$. Feature loss is applied here as well. Note that all losses for future time steps ($t \geq 1$) are averaged and then added to the loss for the current time step ($t = 0$), since the action executed immediately should be our key target to optimize.

To help the agent better estimate its current state, we add a speed prediction head to predict current speed $s$ from the image feature and a value prediction head to predict the expected return estimated by the expert, similarly as in [60]. We take the $L_1$ loss for the speed prediction and $L_2$ loss for the value prediction, denoting their weighted sum as $\mathcal{L}_{aux}$.

The overall loss is as follows, as weighted by $\lambda_{traj}, \lambda_{ctl}, \lambda_{aux}$:

$$\mathcal{L} = \lambda_{traj} \cdot \mathcal{L}_{traj} + \lambda_{ctl} \cdot \mathcal{L}_{ctl} + \lambda_{aux} \cdot \mathcal{L}_{aux}. \tag{6}$$

### 3.4 Output Fusion

We have two forms of output representations from our TCP framework: the planned trajectory and the predicted control. To further combine their advantages, we devise a situation-based fusion strategy as depicted in Algorithm 1. Specifically, denote $\alpha$ as a combination weight whose value is between 0 to 0.5, in a certain situation where one representation is more suitable according to our prior belief, we combine the results from trajectory and control predictions by taking average with weight $\alpha$ so that the more suitable one takes up more weight $(1 - \alpha)$. Note that the combination weight $\alpha$ indeed does not need to be

**Algorithm 1:** Situation based fusion scheme to combine the two output paradigms

**Input:** sensory input $\mathbf{i}$, speed of the ego vehicle $v$, high level navigation information $\mathbf{g}$.
**Hyper parameters:** combination weight $\alpha \in [0, 0.5]$
**Output:** final control signals $\mathbf{a}$

$\{\mathbf{wp}_t\}_{t=0}^{K}, \mathbf{a}^{\text{ctl}} \leftarrow \text{TCP}(\mathbf{i}, v, \mathbf{g})$
$\mathbf{a}^{\text{traj}} \leftarrow \text{Low-level Controller} (\{\mathbf{wp}_t\}_{t=0}^{K})$
Get current *situation*
**if** *situation* is *trajectory specialized* **then**
  | $\mathbf{a} \leftarrow \alpha \times \mathbf{a}^{\text{ctl}} + (1 - \alpha) \times \mathbf{a}^{\text{traj}}$
**else**
  | $\mathbf{a} \leftarrow \alpha \times \mathbf{a}^{\text{traj}} + (1 - \alpha) \times \mathbf{a}^{\text{ctl}}$
**end**

a constant or symmetric, which means we can set it to different values under different situations or different for specific control signals. In our experiment, we choose the *situation* according to whether the ego vehicle is turning, implying that if it is turning, the *situation* is *control specialized* otherwise *trajectory specialized*.

## 4 Experiments

### 4.1 Experimental Setup

**Task & Evaluation metrics.** Our method is validated and tested in the CARLA driving simulator [20]. Given a route defined by a sequence of sparse navigation points together with high level commands (straight, turn left/right, lane changing, and lane following), the closed-loop driving task

Table 1: Evaluation on the public CARLA Leaderboard [1] (accessed in May 2022). Our method TCP and TCP-Ens achieve a driving score of 69.714 and 75.137 respectively with only a monocular camera. More detailed infraction statistics can be found in the Supplementary.

| Rank | Method | Sensor Inputs | | Key Metrics ↑ | | |
|------|--------|---------------|-------|---------------|------------------|-------------------|
| | | #Cameras | LiDAR | Driving Score | Route Completion | Infraction Score |
| 1 | **TCP-Ens** (ours) | 1 | ✗ | **75.137** | 85.629 | **0.873** |
| 1 | **TCP** (ours) | 1 | ✗ | **69.714** | 82.962 | **0.851** |
| 1 | **TCP-SB** (ours) | 1 | ✗ | **68.695** | 82.957 | **0.833** |
| 2 | LAV [10] | 4 | ✓ | 61.846 | **94.459** | 0.640 |
| 3 | Transfuser | 3 | ✓ | 61.181 | 86.694 | 0.714 |
| 4 | Latent Transfuser | 3 | ✗ | 45.029 | 75.366 | 0.618 |
| 5 | GRIAD [9] | 3 | ✗ | 36.787 | 61.855 | 0.597 |
| 6 | Transfuser+ [29] | 4 | ✓ | 34.577 | 69.841 | 0.562 |
| 7 | WoR [12] | 4 | ✗ | 31.370 | 57.647 | 0.557 |
| 8 | MaRLn [50] | 1 | ✗ | 24.980 | 46.968 | 0.518 |
| 9 | NEAT [15] | 3 | ✗ | 21.832 | 41.707 | 0.650 |

Table 2: Comparison between the control and trajectory only model in terms of infractions frequency. TurnRatio means the corresponding ratio of happening during turning.

| Model | Driving Score | Collisions vehicles | | Collisions layout | | Off-road infractions | | Agent blocked | |
|-------|---------------|---------------------|----------|-------------------|----------|----------------------|----------|---------------|----------|
| | | #/km ↓ | TurnRatio | #/km ↓ | TurnRatio | #/km ↓ | TurnRatio | #/km ↓ | TurnRatio |
| Control-Only | 32.45±2.23 | 1.25 | 50.90% | **0.23** | 10.00% | **0.59** | 46.15% | **0.41** | 50.00% |
| Trajectory-Only | 28.29±3.03 | **0.85** | 38.70% | 0.77 | 64.20% | 0.74 | 62.90% | 0.77 | 64.20% |

requires the autonomous agent to drive towards the destination point. It is designed to simulate realistic traffic situations and includes different challenging scenarios such as obstacle avoidance, crossing an unsignalized intersection, and sudden control loss. There are three major metrics: Driving Score, Route Completion, and Infraction Score. Route Completion is the percentage of the route completed by the autonomous agent. Infraction Score measures the number of infractions made along the route, with pedestrians, vehicles, road layouts, red lights, and *etc*. Driving Score is the main metric which is the product of Route Completion and Infraction Score.

**Dataset.** We use randomly generated routes under random weather conditions to collect 420K data in the 8 public towns offered by the CARLA simulator. Similar to [10], we train TCP on 189K of data in 4 out of 8 towns (Town01, Town03, Town04, and Town06) for ablations and train with all 420K data for our online leaderboard submission.

## 4.2 State-of-the-art Comparison

Table 1 shows the result of the comparison between our method and the top 8 entries on the public CARLA Leaderboard [1]. We report the results of TCP and two variants. **TCP-SB** replaces shared encoders of TCP with two separate ones for two branches, and **TCP-Ens** is the ensemble of TCP and TCP-SB. Our method TCP-Ens ranks first on the leaderboard with a 75.137 driving score and highest infraction score, and TCP alone also surpasses prior methods. Note that our method only uses a monocular camera while the top 2-4 methods all use multiple cameras and a LiDAR. Our driving score is 50.157 higher than the second-best monocular camera method, MaRLn [50]. Our route completion is slightly inferior to the LiDAR candidates - one reason is that methods using LiDAR may have a better object detection ability. Based on the detection results, they usually adopt a crawling strategy, indicating that the vehicle would move slowly when it has stopped for a long time and there are no obstacles ahead. As described in [29], this could alleviate ego vehicle's blocking problems to boost the route completion performance.

Figure 4: The trajectory-guided attention maps in two cases. In each case (row), from the left to right we show that the input image with the predicted trajectory (the first waypoint is projected out of the image), the predicted trajectory in the top-down view, the attention map $\mathbf{w}_1$, the attention map $\mathbf{w}_3$.

Table 3: Ablative study on the effectiveness of different components design of our model.

| Exp. | Driving Score | Route Completion | Infraction Score |
|---|---|---|---|
| Control | 32.45±2.23 | 76.54±3.22 | 0.45±0.03 |
| + traj-task | 34.98±1.96 | 81.32±5.50 | 0.49±0.05 |
| + temporal | 42.87±4.77 | **87.51**±3.63 | 0.49±0.07 |
| + traj-attn | 46.08±3.47 | 84.95±1.84 | 0.56±0.03 |
| + fusion | **57.01**±1.88 | 85.27±1.20 | **0.67**±0.01 |

Table 4: Comparison between MTL and ensemble methods ($\alpha$ is 0.3 for all experiments).

| Exp. | Driving Score | #Param. | FLOPs | FPS |
|---|---|---|---|---|
| Ensemble | 45.03±1.28 | 46.81M | 17.07G | 69.47 |
| MTL | 48.27±0.58 | 23.58M | 8.54G | 133.30 |
| TCP-SB | 52.46±4.66 | 47.26M | 17.07G | 69.35 |
| TCP | 57.01±1.88 | 25.77M | 8.54G | 125.71 |
| TCP-Ens | 59.09±3.66 | 73.03M | 25.61G | 44.70 |

## 4.3 Control vs. Trajectory

In this section, we conduct quantitative experiments to compare the **Control-Only** model and the **Trajectory-Only** model to demonstrate their advantages and disadvantages. For both models, we use the same setting except for the output head and its corresponding loss. We use a ResNet-34 to encode visual inputs and a measurement module to encode the navigation information. Similar to [60], we add speed and value heads as auxiliary tasks to help the model better encode the environment. For **Control-Only**, we predict the control distribution based on the concatenated latent feature from the two encoders. As for **Trajectory-Only**, we feed the feature to a GRU decoder to generate waypoints. As shown in Table 2, though Trajectory-Only collides with vehicles less frequently than Control-Only, it has more layout collisions, off-road infractions, and agent blocks. We also count the ratio of each kind of infraction that occurs during turning. It can be observed that for Trajectory-Only, a large portion of such infractions happen when the ego agent is turning compared to Control-Only. This has verified that Trajectory-Only performs worse when the agent is turning, which is probably caused by the unsatisfactory trajectory following performance of simple PID controllers as discussed in Sec. 1. As for the fact that Control-Only has a higher vehicle-collision rate, it is because the model focuses on the current time step and the reaction to potential collisions tends to be late, as depicted in Sec. 1 as well. The results above further validate the necessity of combining the two output paradigms.

## 4.4 Ablative Study and Visualization

**Component analysis.** We first validate the effectiveness of the trajectory-guided multi-step control prediction design, as shown in Table 3. We only employ the control branch output except for the last complete one when fusion is applied for these ablations. Adding a trajectory branch as an auxiliary task improves the performance by 2.5 points. The multi-step predictions with our temporal module greatly help with 7.9 points gain, and adding the trajectory-guided attention further acquires an improvement of 3.2 points. Finally, applying our situation based fusion scheme ($\alpha$ is set to 0.3) significantly boosts the infraction score, leading the overall driving score to 57.

**Multi-task vs. Ensemble.** The comparison regarding their performances and computational complexity is given in Table 4. **Ensemble** denotes directly combining the outputs of Control-Only and Trajectory-Only with our situation based fusion scheme. **MTL** represents the model with a shared CNN backbone and measurement encoders followed by a trajectory branch and a control branch, but the control branch predicts current step prediction only and there are no interactions between the two branches. We conclude that directly combining two models with our fusion scheme greatly improves the performance, and using an MTL approach works better than ensemble but with a much

smaller model size and GFLOPs. A conventional ensemble approach to combine results from TCP and TCP-SB as TCP-Ens brings further performance gain at the cost of computational complexity.

**Situation based fusion weight.** We investigate the choice of the combination weight $\alpha$ in the situation based fusion scheme and show the box plot of the driving scores in the figure to the right. Besides $\alpha \in [0, 0.5]$, we additionally test 0.7 and 1, meaning that two results are conversely mixed with our *specialization* definition. We see that only using the control from the specialized branch $\alpha = 0$ performs poorly while directly taking the average or fusing conversely still has comparable results. One reason is that the *situation* criterion used here is whether the vehicle is turning, making most cases *trajectory specialized*, and the stronger control branch is not utilized enough if $\alpha$ is small. Note that the situation based fusion scheme is general and flexible, and the criterion or $\alpha$ value used here is relatively coarse.
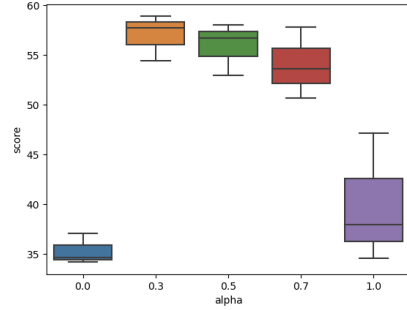


Figure 5: Box plot of the driving score with different $\alpha$ values (3 trials for each $\alpha$).

**Visualization.** Fig. 4 visualizes the trajectory-guided attention maps. The trajectory branch provides location-related information to guide the control branch to focus on important regions which are useful for future control prediction. See more qualitative results in the Supplementary.

## 5 Conclusion

In this work, we study two learning and prediction paradigms based on trajectory and direct control, respectively, for end-to-end autonomous driving. We propose a unified framework comprised of a trajectory branch and a novel multi-step control branch with interactions in between. We design a situation based fusion scheme to combine the results from two branches. Our method with only a monocular camera has achieved state-of-the-art performance on the CARLA Leaderboard.

## Acknowledgments

## References

[1] CARLA autonomous driving leaderboard. https://leaderboard.carla.org/, 2022. 3, 8

[2] Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. Multi-task feature learning. In *NeurIPS*, 2006. 2, 3

[3] Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016. 3

[4] Mariusz Bojarski, Chenyi Chen, Joyjit Daw, Alperen Değirmenci, Joya Deri, Bernhard Firner, Beat Flepp, Sachin Gogri, Jesse Hong, Lawrence Jackel, et al. The nvidia pilotnet experiments. *arXiv preprint arXiv:2010.08776*, 2020. 1

[5] Eduardo F Camacho and Carlos Bordons Alba. *Model predictive control*. Springer science & business media, 2013. 2

[6] Manuel Carranza-García, Pedro Lara-Benítez, Jorge García-Gutiérrez, and José C Riquelme. Enhancing object detection for autonomous driving by optimizing anchor generation and addressing class imbalance. *Neurocomputing*, 2021. 4

[7] Rich Caruana. Multitask learning. *Machine learning*, 1997. 2, 3

[8] Sergio Casas, Abbas Sadat, and Raquel Urtasun. Mp3: A unified model to map, perceive, predict and plan. In *CVPR*, 2021. 2, 3

[9] Raphael Chekroun, Marin Toromanoff, Sascha Hornauer, and Fabien Moutarde. Gri: General reinforced imitation and its application to vision-based autonomous driving. *arXiv preprint arXiv:2111.08575*, 2021. 1, 3, 8

[10] Dian Chen and Philipp Krähenbühl. Learning from all vehicles. In *CVPR*, 2022. 1, 3, 4, 8

[11] Dian Chen, Brady Zhou, Vladlen Koltun, and Philipp Krähenbühl. Learning by cheating. In *CoRL*, 2020. 1, 2, 3

[12] Dian Chen, Vladlen Koltun, and Philipp Krähenbühl. Learning to drive from a world on rails. In *ICCV*, 2021. 1, 3, 8

[13] Li Chen, Chonghao Sima, Yang Li, Zehan Zheng, Jiajie Xu, Xiangwei Geng, Hongyang Li, Conghui He, Jianping Shi, Yu Qiao, et al. Persformer: 3d lane detection via perspective transformer and the openlane benchmark. *arXiv preprint arXiv:2203.11089*, 2022. 2, 4

[14] Sumanth Chennupati, Ganesh Sistu, Senthil Yogamani, and Samir A Rawashdeh. Multinet++: Multi-stream feature aggregation and geometric loss strategy for multi-task learning. In *CVPRW*, 2019. 4

[15] Kashyap Chitta, Aditya Prakash, and Andreas Geiger. Neat: Neural attention fields for end-to-end autonomous driving. In *ICCV*, 2021. 1, 2, 3, 8

[16] Kyunghyun Cho, Bart van Merrienboer, Çaglar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *EMNLP*, 2014. 2, 5

[17] Felipe Codevilla, Matthias Müller, Antonio López, Vladlen Koltun, and Alexey Dosovitskiy. End-to-end driving via conditional imitation learning. In *ICRA*, 2018. 1, 2, 3, 4

[18] Felipe Codevilla, Eder Santana, Antonio M López, and Adrien Gaidon. Exploring the limitations of behavior cloning for autonomous driving. In *ICCV*, 2019. 1, 2, 3

[19] Thomas G Dietterich. Ensemble methods in machine learning. In *MCS*, 2000. 4

[20] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *CoRL*, 2017. 3, 7

[21] Hongyan Guo, Dongpu Cao, Hong Chen, Zhenping Sun, and Yunfeng Hu. Model predictive path following control for autonomous cars considering a measurable disturbance: Implementation, testing, and verification. *MSSP*, 2019. 2

[22] Jeffrey Hawke, Richard Shen, Corina Gurau, Siddharth Sharma, Daniele Reda, Nikolay Nikolov, Przemysław Mazur, Sean Micklethwaite, Nicolas Griffiths, Amar Shah, et al. Urban driving with conditional imitation learning. In *ICRA*, 2020. 3

[23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5

[24] Yuenan Hou, Zheng Ma, Chunxiao Liu, and Chen Change Loy. Learning to steer by mimicking features from heterogeneous auxiliary networks. In *AAAI*, 2019. 3, 4

[25] Shengchao Hu, Li Chen, Penghao Wu, Hongyang Li, Junchi Yan, and Dacheng Tao. St-p3: End-to-end vision-based autonomous driving via spatial-temporal feature learning. In *ECCV*, 2022. 2, 3

[26] Sebastian Huch, Aybike Ongel, Johannes Betz, and Markus Lienkamp. Multi-task end-to-end self-driving architecture for cav platoons. *Sensors*, 2021. 3

[27] Keishi Ishihara, Anssi Kanervisto, Jun Miura, and Ville Hautamaki. Multi-task learning with attention for end-to-end autonomous driving. In *CVPR*, 2021. 3, 4

[28] Robert A Jacobs, Michael I Jordan, Steven J Nowlan, and Geoffrey E Hinton. Adaptive mixtures of local experts. *Neural computation*, 1991. 4

[29] Bernhard Jaeger. Expert drivers for autonomous driving. Master's thesis, University of Tübingen, 2021. 1, 2, 3, 4, 8

[30] Alex Kendall, Jeffrey Hawke, David Janz, Przemyslaw Mazur, Daniele Reda, John-Mark Allen, Vinh-Dieu Lam, Alex Bewley, and Amar Shah. Learning to drive in a day. In *ICRA*, 2019. 3

[31] Inhan Kim, Hyemin Lee, Joonyeong Lee, Eunseop Lee, and Daijin Kim. Multi-task learning with future states for vision-based autonomous driving. In *ACCV*, 2020. 3, 4

[32] Inhan Kim, Joonyeong Lee, and Daijin Kim. Learning mixture of domain-specific experts via disentangled factors for autonomous driving. In *AAAI*, 2022. 4

[33] Philipp Krähenbühl and Vladlen Koltun. Learning to propose objects. In *CVPR*, 2015. 4

[34] Varun Ravi Kumar, Senthil Yogamani, Hazem Rashed, Ganesh Sitsu, Christian Witt, Isabelle Leang, Stefan Milz, and Patrick Mäder. Omnidet: Surround view cameras based multi-task visual perception network for autonomous driving. *RA-L*, 2021. 2, 4

[35] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *NeurIPS*, 2017. 4

[36] Yingwei Li, Adams Wei Yu, Tianjian Meng, Ben Caine, Jiquan Ngiam, Daiyi Peng, Junyang Shen, Bo Wu, Yifeng Lu, Denny Zhou, et al. Deepfusion: Lidar-camera deep fusion for multi-modal 3d object detection. In *CVPR*, 2022. 4

[37] Zhihao Li, Toshiyuki Motoyoshi, Kazuma Sasaki, Tetsuya Ogata, and Shigeki Sugano. Rethinking self-driving: Multi-task knowledge for better generalization and accident explanation ability. *arXiv preprint arXiv:1809.11100*, 2018. 3

[38] Ming Liang, Bin Yang, Yun Chen, Rui Hu, and Raquel Urtasun. Multi-task multi-sensor fusion for 3d object detection. In *CVPR*, 2019. 2, 4

[39] Xiaodan Liang, Tairui Wang, Luona Yang, and Eric Xing. Cirl: Controllable imitative reinforcement learning for vision-based self-driving. In *ECCV*, 2018. 1, 3

[40] Urs Muller, Jan Ben, Eric Cosatto, Beat Flepp, and Yann Cun. Off-road obstacle avoidance through end-to-end learning. In *NeurIPS*, 2005. 3

[41] Urs Muller, Jan Ben, Eric Cosatto, Beat Flepp, and Yann Cun. Off-road obstacle avoidance through end-to-end learning. In *NeurIPS*, 2005. 2

[42] Eshed Ohn-Bar, Aditya Prakash, Aseem Behl, Kashyap Chitta, and Andreas Geiger. Learning situational driving. In *CVPR*, 2020. 1, 4

[43] Ashwini Pokle, Roberto Martín-Martín, Patrick Goebel, Vincent Chow, Hans M Ewald, Junwei Yang, Zhenkai Wang, Amir Sadeghian, Dorsa Sadigh, Silvio Savarese, et al. Deep local trajectory replanning and control for robot navigation. In *ICRA*, 2019. 3

[44] Dean A Pomerleau. Alvinn: An autonomous land vehicle in a neural network. In *NeurIPS*, 1988. 2, 3

[45] Aditya Prakash, Aseem Behl, Eshed Ohn-Bar, Kashyap Chitta, and Andreas Geiger. Exploring data aggregation in policy learning for vision-based urban autonomous driving. In *CVPR*, 2020. 3, 4

[46] Aditya Prakash, Kashyap Chitta, and Andreas Geiger. Multi-modal fusion transformer for end-to-end autonomous driving. In *CVPR*, 2021. 1, 3, 5

[47] Georgia Rajamanoharan, Aytaç Kanacı, Minxian Li, Shaogang Gong, et al. Multi-task mutual learning for vehicle re-identification. In *CVPR*, 2019. 4

[48] Nicholas Rhinehart, Rowan McAllister, and Sergey Levine. Deep imitative models for flexible inference, planning, and control. In *ICLR*, 2019. 1

[49] Abbas Sadat, Sergio Casas, Mengye Ren, Xinyu Wu, Pranaab Dhawan, and Raquel Urtasun. Perceive, predict, and plan: Safe motion planning through interpretable semantic representations. In *ECCV*, 2020. 3

[50] Marin Toromanoff, Emilie Wirbel, and Fabien Moutarde. End-to-end model-free reinforcement learning for urban driving using implicit affordances. In *CVPR*, 2020. 3, 8

[51] Apoorv Vyas, Nataraj Jammalamadaka, Xia Zhu, Dipankar Das, Bharat Kaul, and Theodore L Willke. Out-of-distribution detection using an ensemble of self supervised leave-out classifiers. In *ECCV*, pages 550–564, 2018. 4

[52] Dong Wu, Manwen Liao, Weitian Zhang, and Xinggang Wang. Yolop: You only look once for panoptic driving perception. *arXiv preprint arXiv:2108.11250*, 2021. 4

[53] Yi Xiao, Felipe Codevilla, Akhil Gurram, Onay Urfalioglu, and Antonio M López. Multimodal end-to-end autonomous driving. *T-ITS*, 2020. 3

[54] Huazhe Xu, Yang Gao, Fisher Yu, and Trevor Darrell. End-to-end learning of driving models from large-scale video datasets. In *CVPR*, 2017. 3

[55] Jie Xu, Wei Wang, Hanyuan Wang, and Jinhong Guo. Multi-model ensemble with rich spatial information for object detection. *Pattern Recognition*, 2020. 4

[56] Zhengyuan Yang, Yixuan Zhang, Jerry Yu, Junjie Cai, and Jiebo Luo. End-to-end multi-modal multi-task vehicle control for self-driving cars with visual perceptions. In *ICPR*, 2018. 3, 4

[57] Éloi Zablocki, Hédi Ben-Younes, Patrick Pérez, and Matthieu Cord. Explainability of vision-based autonomous driving systems: Review and challenges. *arXiv preprint arXiv:2101.05307*, 2021. 2

[58] Wenyuan Zeng, Wenjie Luo, Simon Suo, Abbas Sadat, Bin Yang, Sergio Casas, and Raquel Urtasun. End-to-end interpretable neural motion planner. In *CVPR*, 2019. 3

[59] Jimuyang Zhang and Eshed Ohn-Bar. Learning by watching. In *CVPR*, 2021. 1, 3

[60] Zhejun Zhang, Alexander Liniger, Dengxin Dai, Fisher Yu, and Luc Van Gool. End-to-end urban driving by imitating a reinforcement learning coach. In *ICCV*, 2021. 1, 3, 4, 7, 9

[61] Albert Zhao, Tong He, Yitao Liang, Haibin Huang, Guy Van den Broeck, and Stefano Soatto. Sam: Squeeze-and-mimic networks for conditional visual driving policy learning. In *CoRL*, 2021. 3

[62] Yinuo Zhao, Kun Wu, Zhiyuan Xu, Zhengping Che, Qi Lu, Jian Tang, and Chi Harold Liu. Cadre: A cascade deep reinforcement learning framework for vision-based autonomous urban driving. In *AAAI*, 2022. 3

[63] Zeyu Zhu and Huijing Zhao. Multi-task conditional imitation learning for autonomous navigation at crowded intersections. *arXiv preprint arXiv:2202.10124*, 2022. 3, 4

# Supplementary Materials:
# Trajectory-guided Control Prediction for End-to-end Autonomous Driving: A Simple yet Strong Baseline

**Penghao Wu**[*][†]
Shanghai AI Laboratory
Shanghai Jiao Tong University
wupenghaocraig@sjtu.edu.cn

**Xiaosong Jia**[*]
Shanghai AI Laboratory
Shanghai Jiao Tong University
jiaxiaosong@sjtu.edu.cn

**Li Chen**[*]
Shanghai AI Laboratory
lichen@pjlab.org.cn

**Junchi Yan**[†]
Shanghai Jiao Tong University
yanjunchi@sjtu.edu.cn

**Hongyang Li**
Shanghai AI Laboratory
Shanghai Jiao Tong University
lihongyang@pjlab.org.cn

**Yu Qiao**
Shanghai AI Laboratory
qiaoyu@pjlab.org.cn

In this Supplementary document, we first provide a detailed description of the dataset in Sec. A. Implementation and training details are in Sec. B. We show detailed infraction statistics for both leaderboard results and ablation studies, and qualitative results in Sec. C. Last, we discuss limitations, common failure cases, and possible future directions and potential social impact of our work in Sec. D.

## A  Dataset

### A.1  Dataset Collection

We use CARLA 0.9.10.1 for data collection and testing. We use Roach [16] as the expert to collect data. In order to improve the obstacle avoidance ability of the expert, we additionally add a rule-based vehicle and pedestrian detector adopted from Transfuser [12] to avoid possible collisions. Each route is generated randomly with length ranging from 50 meters to 300 meters. We use the scenario configurations provided in [12]. We terminate each route if the expert makes a collision or runs a red light. Last few frames for such routes are discarded. The data samples are stored at 2HZ.

### A.2  Dataset Statistics

Detailed statistics for each town and their descriptions are provided in Table 1. As stated in the main paper, we train on all eight towns for the leaderboard submission. For our ablation experiments, we train on four towns (Town01, Town03, Town04, and Town06) and test on the designed four routes with four different weathers in Town02 and Town05, as does in [3].

## B  Implementation Details

We use ResNet-34 [8] pretrained on ImageNet [6] as the image encoder. The size of the input image is $900 \times 256$ and the FOV of the camera is set as $100°$. We choose $K$ being 4 meaning four future steps at 2HZ are predicted for both the trajectory branch and the control branch. Detailed network structure is presented in Table 5. We follow the same PID setting as [5], where the PID parameters are exquisitely tuned, *i.e.*, $K_p = 5.0$, $K_i = 0.5$, $K_d = 1.0$ for the longitudinal PID controller and $K_p = 0.75$, $K_i = 0.75$, $K_d = 0.3$ for the lateral PID controller. The weights for different loss terms are as follows: $\lambda_F = 0.05$, $\lambda_{traj} = 1$, $\lambda_{ctl} = 1$, $\lambda_{aux} = 0.05$, and 0.001 for speed and value regression

Table 1: Detailed statistics of the number of samples, the number of dynamic agents added, and a brief description of each town.

| Town Name | #Samples | #Dynamic Agents | Description |
|---|---|---|---|
| Town01 | 50384 | 120 | a basic town with T junction |
| Town02 | 55943 | 100 | similar to Town01 but smaller |
| Town03 | 42771 | 120 | a complex town |
| Town04 | 47954 | 200 | a highway loop and a small town |
| Town05 | 53684 | 120 | a squared-grid town with multiple lanes |
| Town06 | 48415 | 150 | long highways |
| Town07 | 51549 | 110 | a rural enviroment with narrow roads |
| Town10 | 59898 | 120 | a city with various environments |

respectively. For all experiments, we train TCP on 4 GeForce RTX 3090 GPUs. We use the Adam optimizer [10] with a learning rate of $1 \times 10^{-4}$ and weight decay of $1 \times 10^{-7}$ for all experiments. We train all models with batch size 128 for 60 epochs, and the learning rate is reduced by a factor of 2 after 30 epochs.

In the situation based fusion scheme, we choose whether the vehicle is turning as the criterion of the *situation*. Specifically, we calculate the absolute values of steer actions within the past 1 second. If half of them are larger than 0.1, we assume the vehicle is turning so the *situation* is *control specialized*, otherwise *trajectory specialized*. For the online CARLA Leaderboard [1] submission, we use an asymmetric fusion scheme. If the *situation* is *trajectory specialized*, we set $\alpha = 0.5$, and $\alpha = 0$ when it is *control specialized*. We take the maximum of the *brake* control instead of taking the average. For the ensemble submission TCP-Ens, we also take the maximum of *brake* value from different models and take the average for *steer* and *throttle*.

## C  Experiments

### C.1  Validation Protocol Details

We use the same validation routes as LAV [3]. This includes 4 routes in total, 2 from Town02 and 05 each. Each route is tested under 4 different weathers (ClearNoon, CloudySunset, SoftRainDawn, HardRainNight) and is repeated for 3 times, resulting in 48 routes in total. Random scenarios are added from the official CARLA leaderboard repo (all_towns_traffic_scenarios_public.json). The time-limit for agent blocking is reduced from 300 seconds to 60 seconds to save time.

### C.2  Detailed Infractions Statistics

In this part, we report detailed infraction statistics of the methods on CARLA Leaderboard in Table 2, and statistics of our ablation experiments in Table 3 and Table 4.

### C.3  Qualitative Results

We show cases of our method performing well in different challenging scenarios in Fig. 1. In the first case, the autonomous agent successfully reacts to the changing of the traffic light in time at the crossing. In the second case, a cyclist suddenly runs across the road right after the ego vehicle has made a right turn, and our agent makes an emergency brake in time, avoiding a collision. In the third case, the ego vehicle is making a right turn while there are other vehicles crossing. It stops and waits for the crossing vehicles to pass and then continues to make the turn. In the last case, our agent is performing an unprotected left turn with oncoming traffic, and it successfully negotiates with the oncoming vehicle.

More visualization examples of the trajectory-guided attention maps are provided in Fig. 2. We also show the GradCam [13] and EigenCam [11] visualization of two examples for Control-Only model with multi-step prediction scheme in Fig. 3. For the GradCam visualization, we set the target (which is needed to be maximized during the calculation of GradCam) be the negative action loss for the current and future action prediction. Note that GradCam visualizes the regions of the input image that are **important** for predictions by calculating gradients to maximize the target. It **does not** indicate that the model does focus or well capture the region highlighted by GradCam. As shown in Fig. 3, the

Table 2: Detailed statistics of the evaluation on the public CARLA Leaderboard [1] (accessed in May 2022). Driving Score, Route Completion, and Infraction Penalty are higher the better. For other metrics, lower values are desired. The collisions, infractions, and agent blocked related metrics are given as the number of events per kilometer. Our method outperforms other methods by a large margin in terms of Driving Score and Route Completion. We also have the best scores for metrics of collisions vehicle, collisions pedestrian, collisions layout, and off-road infractions among all methods.

| Rank | Method | Driving Score | Route Completion | Infraction Penalty | Collisions Vehicle | Collisions Pedestrian | Collisions Layout | Red light Infractions | Off-road Infractions | Agent Blocked |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | **TCP-Ens** (ours) | **75.137** | 85.629 | **0.873** | 0.316 | **0.000** | **0.000** | 0.089 | 0.038 | 0.537 |
| 1 | **TCP** (ours) | 69.714 | 82.962 | 0.851 | **0.220** | 0.006 | 0.034 | 0.083 | **0.017** | 0.564 |
| 1 | **TCP-SB** (ours) | 68.695 | 82.957 | 0.833 | 0.250 | **0.000** | 0.111 | 0.066 | 0.026 | 0.528 |
| 2 | LAV [3] | 61.846 | **94.459** | 0.640 | 0.696 | 0.038 | 0.017 | 0.166 | 0.252 | **0.104** |
| 3 | Transfuser | 61.181 | 86.694 | 0.714 | 0.814 | 0.036 | 0.007 | **0.046** | 0.228 | 0.428 |
| 4 | Latent Transfuser | 45.029 | 75.366 | 0.618 | 1.259 | 0.034 | 0.098 | 0.102 | 0.288 | 0.757 |
| 5 | GRIAD [2] | 36.787 | 61.855 | 0.597 | 2.772 | **0.000** | 0.407 | 0.484 | 1.388 | 0.842 |
| 6 | Transfuser+ [9] | 34.577 | 69.841 | 0.562 | 0.703 | 0.045 | 0.025 | 0.750 | 0.185 | 2.406 |
| 7 | WoR [4] | 31.370 | 57.647 | 0.557 | 1.346 | 0.606 | 1.017 | 0.791 | 0.963 | 0.473 |
| 8 | MaRLn [14] | 24.980 | 46.968 | 0.518 | 2.329 | **0.000** | 2.472 | 0.550 | 1.823 | 0.936 |
| 9 | NEAT [5] | 21.832 | 41.707 | 0.650 | 0.742 | 0.042 | 0.617 | 0.700 | 2.680 | 5.225 |

Table 3: Detailed infraction statistics of the ablation on the effectiveness of the trajectory-guided multi-step control prediction design.

| Exp. | Driving Score | Route Completion | Infraction Penalty | Collisions Vehicle | Collisions Pedestrian | Collisions Layout | Red light Infractions | Off-road Infraction | Agent Blocked |
|---|---|---|---|---|---|---|---|---|---|
| Control | 32.45±2.23 | 76.54±3.22 | 0.45±0.03 | 1.24±0.06 | 0.00±0.00 | 0.23±0.09 | 0.18±0.05 | 0.59±0.06 | 0.41±0.11 |
| + traj-task | 34.98±1.96 | 81.32±5.50 | 0.49±0.05 | 1.39±0.15 | 0.00±0.00 | 0.15±0.07 | 0.11±0.04 | 0.39±0.04 | 0.38±0.10 |
| + temporal | 42.87±4.77 | **87.51**±3.63 | 0.49±0.07 | 1.14±0.25 | 0.00±0.00 | 0.20±0.07 | 0.18±0.04 | 0.18±0.05 | 0.22±0.03 |
| + traj-attn | 46.08±3.47 | 84.95±1.84 | 0.56±0.03 | 0.90±0.20 | 0.00±0.00 | **0.04**±0.06 | 0.14±0.07 | 0.54±0.06 | 0.29±0.08 |
| + fusion | **57.01**±1.88 | 85.27±1.20 | **0.67**±0.01 | **0.37**±0.10 | 0.00±0.00 | 0.08±0.03 | **0.10**±0.03 | **0.14**±0.06 | **0.20**±0.03 |

Table 4: Detailed infraction statistics of the experiments of the comparison between MTL and ensemble methods.

| Exp. | Driving Score | Route Completion | Infraction Penalty | Collisions Vehicle | Collisions Pedestrian | Collisions Layout | Red light Infractions | Off-road Infraction | Agent Blocked |
|---|---|---|---|---|---|---|---|---|---|
| Ensemble | 45.03±1.28 | 79.30±5.13 | 0.59±0.04 | 0.62±0.09 | 0.00±0.00 | 0.22±0.03 | 0.22±0.07 | 0.28±0.03 | 0.35±0.10 |
| MTL | 48.27±0.58 | 81.62±2.74 | 0.60±0.02 | 0.51±0.07 | 0.00±0.00 | 0.26±0.06 | **0.06**±0.05 | 0.36±0.03 | 0.28±0.09 |
| TCP-SB | 52.46±4.66 | 83.94±3.75 | 0.64±0.04 | 0.53±0.14 | 0.00±0.00 | 0.08±0.03 | 0.13±0.09 | **0.06**±0.00 | 0.29±0.04 |
| TCP | 57.01±1.88 | 85.27±1.20 | 0.67±0.01 | **0.37**±0.10 | 0.00±0.00 | 0.08±0.03 | 0.10±0.03 | 0.14±0.06 | **0.20**±0.03 |
| TCP-Ens | **59.09**±3.66 | **87.02**±2.02 | **0.70**±0.03 | 0.41±0.19 | 0.00±0.00 | **0.00**±0.00 | 0.10±0.13 | 0.18±0.05 | 0.27±0.06 |

GradCam heat-map for current action prediction ($a_0$) focuses on regions close to the current location of the ego vehicle, while the heat-map for future prediction ($a_1$) focuses on regions further. This indicates that predicting future actions does need to focus on further regions.

However, Control-Only model with multi-step action prediction only aggregates the image feature map once by global average. The pooled feature is then used to predict actions of all time steps. Therefore, it is **not** realistic to highlight corresponding important regions for each step. We use the EigenCam to visualize the 2D image feature map. It is a gradient-free visualization method to directly calculate the heat-map by projecting the feature map to eigen-vectors. As shown in the last column in Fig. 3, the highlighted region of the 2D image feature map only spans a single area, which is not informative enough for multi-step predictions. It verifies the necessity of re-aggregating the image information with different highlighted regions for each future step, as what we did in Fig. 2.

(a) Respond to traffic lights changes

(b) Emergency brake after a turning

(c) Right turn with crossing traffic

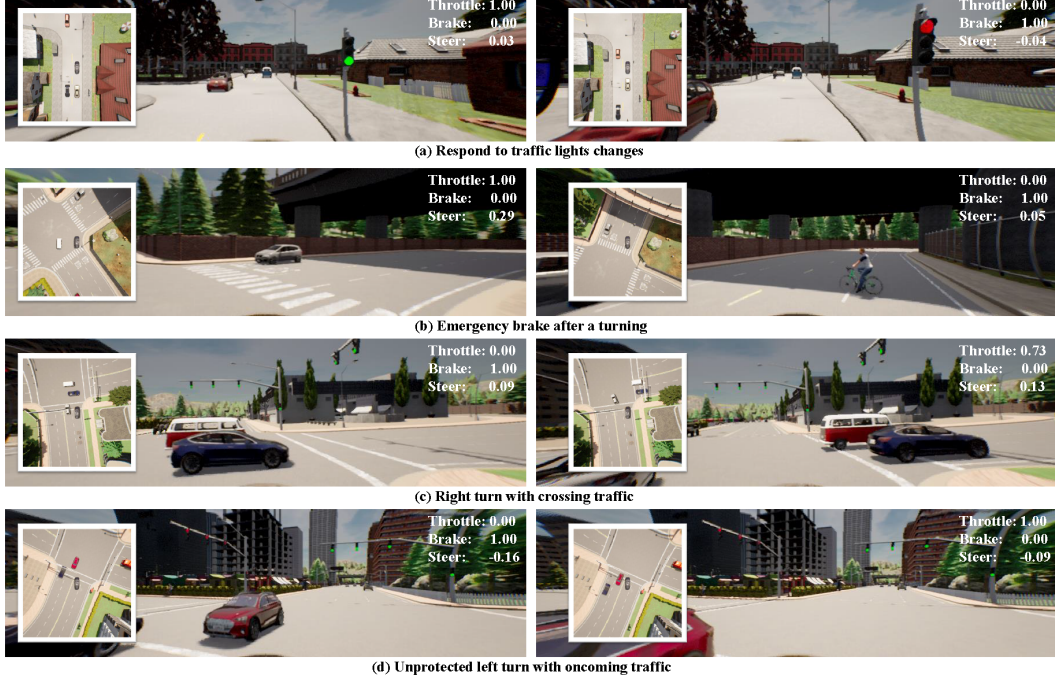(d) Unprotected left turn with oncoming traffic

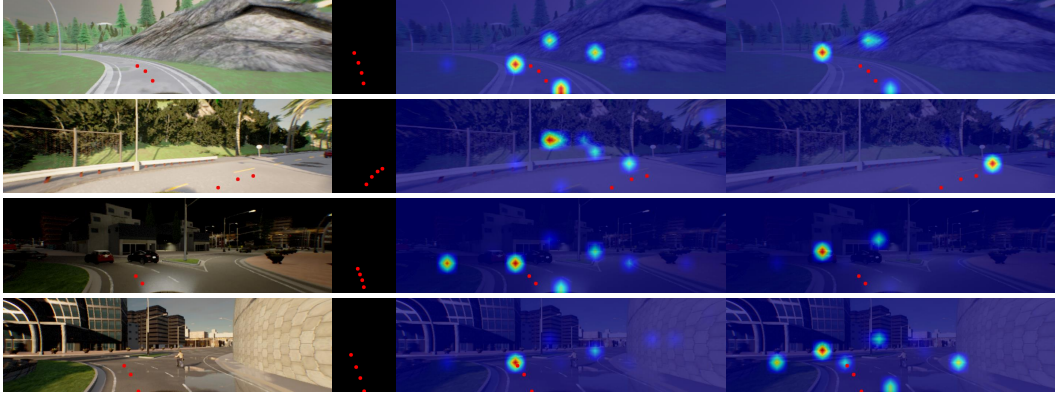Figure 1: Examples of our agent performing well under different challenging scenarios.



Figure 2: More examples of trajectory-guided attention maps. In each case (row), from the left to right we show that the input image with the predicted trajectory (the first waypoint is projected out of the image), the predicted trajectory in the top-down view, the attention map $\mathbf{w}_1$, the attention map $\mathbf{w}_3$.

# D  Discussion

## D.1  Limitations and Future Work

### D.1.1  Failure Cases and Future Work Directions

Our work mainly focuses on combining the two output forms of end-to-end autonomous driving, *i.e.*, trajectory planning and direct control. A detailed and elaborate situation based fusion scheme is based on rules which may require a large number of experiments and specific prior knowledge. A more general or a learning-based adaptive fusion scheme may be a possible future direction.

We also discuss two typical failure cases of TCP in Fig. 4. The first scenario happens when other vehicles initially outside the ego agent's front view rushes into the path with a high speed. It causes

4
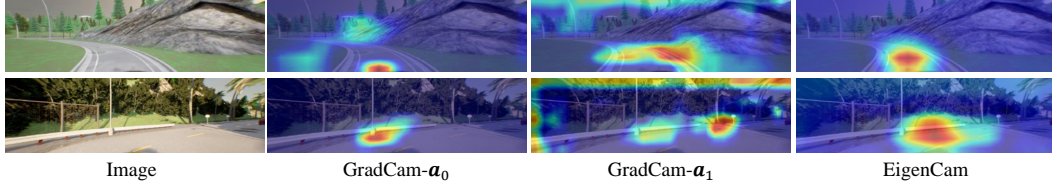
| Image | GradCam-$a_0$ | GradCam-$a_1$ | EigenCam |

Figure 3: Visualization examples of GradCam [13] and EigenCam [11]. From left to right: original image, GradCam heat-map for current action prediction, GradCam heat-map for future action prediction, EigenCam heat-map for the image feature map.



Figure 4: Examples of two failure cases. Top row: the red vehicle runs into the ego path with a high speed, and the ego vehicle fails to take an emergent brake. Down row: The ego agent is waiting for a left turn but occupies part of the opposite lane, causing a block.

a delayed collision when an emergent braking fails. It is because of the limited view of our single camera, hence a straightforward future direction is to add multi-view cameras or a LiDAR input to our agent. Another kind of failures is that the ego agent fails to predict the possible trajectory of other vehicles, resulting in blocking or collisions. Thus explicitly making trajectory predictions of other vehicles and combining it with our trajectory branch is also an interesting direction to further boost the ability of generalization as demonstrated in LAV [3] and LBW [15].

## D.2   Broader Impact

We explore the limitations and advantages of the two conventional output paradigms for end-to-end autonomous driving, present TCP which achieves state-of-the-art performance on the public closed-loop benchmark to push the boundary of the problem. We aim to bring together the two branches of research in this field and provide a unified framework to combine their possible advantages. Our work provides a simple yet effective framework, based on which, new models and techniques can be conveniently integrated and transparently compared. Despite such improvement, we fully understand that our work is by no means perfect and still has many challenges when it comes to real-world application. Our model is trained and tested in the simulator, directly deploying it in the real world will lead to possible traffic accidents which may cause negative societal impacts.

## E   License of Assets

CARLA [7] is an open-source simulator which is under the MIT license and its assets are under the CC-BY license. We integrate part of the official code of Roach [16] which is under the CC-BY-NC 4.0 license into our codebase. The pretrained ResNet model is under the MIT license.

The source code and training data for our work will be publicly available once accepted and they are under the CC-BY-NC 4.0 license.

## References

[1] CARLA autonomous driving leaderboard. https://leaderboard.carla.org/, 2022. 2, 3

[2] Raphael Chekroun, Marin Toromanoff, Sascha Hornauer, and Fabien Moutarde. Gri: General reinforced imitation and its application to vision-based autonomous driving. *arXiv preprint arXiv:2111.08575*, 2021. 3

[3] Dian Chen and Philipp Krähenbühl. Learning from all vehicles. In *CVPR*, 2022. 1, 2, 3, 5

Table 5: Detailed network structure of our TCP model.

| Layer Type | # of Filters | Activation Function | # |
|---|---|---|---|
| Image Encoder | | | |
| ResNet-34 | | | |
| Measurement Encoder | | | |
| FC | 128 | ReLU | × 2 |
| Join_ctrl | | | |
| FC | 512 | ReLU | × 2 |
| FC | 256 | ReLU | × 1 |
| Join_traj | | | |
| FC | 512 | ReLU | × 2 |
| FC | 256 | ReLU | × 1 |
| Speed Head | | | |
| FC | 256 | ReLU | × 2 |
| FC | 1 | ReLU | × 1 |
| Value Head ×2 traj+ctrl | | | |
| FC | 256 | ReLU | × 2 |
| FC | 1 | ReLU | × 1 |
| Temporal Module | | | |
| GRU_cell | hidden size = 256 | | × 1 |
| FC (output) | 256 | ReLU | × 2 |
| Control Policy Head | | | |
| FC | 256 | ReLU | × 2 |
| FC (alpha) | 2 | Softplus | × 1 |
| FC (beta) | 2 | Softplus | × 1 |
| Traj Policy Head | | | |
| GRU_cell | hidden size = 256 | | × 1 |
| FC (output) | 2 | | × 2 |
| Init Att. | | | |
| FC | 256 | ReLU | × 1 |
| FC | 29*8 | Softmax | × 1 |
| Traj Guided Att. | | | |
| FC | 256 | ReLU | × 1 |
| FC | 29*8 | Softmax | × 1 |
| Merge (merge the re-aggregated image feature and hidden state) | | | |
| FC | 512 | ReLU | × 1 |
| FC | 256 | ReLU | × 1 |

[4] Dian Chen, Vladlen Koltun, and Philipp Krähenbühl. Learning to drive from a world on rails. In *ICCV*, 2021. 3

[5] Kashyap Chitta, Aditya Prakash, and Andreas Geiger. Neat: Neural attention fields for end-to-end autonomous driving. In *ICCV*, 2021. 1, 3

[6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 1

[7] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *CoRL*, 2017. 5

[8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1

[9] Bernhard Jaeger. Expert drivers for autonomous driving. Master's thesis, University of Tübingen, 2021. 3

[10] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 2

[11] Mohammed Bany Muhammad and Mohammed Yeasin. Eigen-cam: Class activation map using principal components. In *IJCNN*, 2020. 2, 5

[12] Aditya Prakash, Kashyap Chitta, and Andreas Geiger. Multi-modal fusion transformer for end-to-end autonomous driving. In *CVPR*, 2021. 1

[13] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, 2017. 2, 5

[14] Marin Toromanoff, Emilie Wirbel, and Fabien Moutarde. End-to-end model-free reinforcement learning for urban driving using implicit affordances. In *CVPR*, 2020. 3

[15] Jimuyang Zhang and Eshed Ohn-Bar. Learning by watching. In *CVPR*, 2021. 5

[16] Zhejun Zhang, Alexander Liniger, Dengxin Dai, Fisher Yu, and Luc Van Gool. End-to-end urban driving by imitating a reinforcement learning coach. In *ICCV*, 2021. 1, 5