# Driving in Real Life with Inverse Reinforcement Learning

**Tung Phan-Minh**  **Forbes Howington**  **Ting-Sheng Chu**  **Sang Uk Lee**

**Momchil S. Tomov**  **Nanxiang Li**  **Caglayan Dicle**  **Samuel Findler**

**Francisco Suarez-Ruiz**  **Robert Beaudoin**  **Bo Yang**  **Sammy Omari**

**Eric M. Wolff**
Motional
{tung.phan, forbes.howington, ting-sheng.chu, eric.wolff}@motional.com

## Abstract

In this paper, we introduce the first learning-based planner to drive a car in dense, urban traffic using Inverse Reinforcement Learning (IRL). Our planner, DriveIRL, generates a diverse set of trajectory proposals, filters these trajectories with a lightweight and interpretable safety filter, and then uses a learned model to score each remaining trajectory. The best trajectory is then tracked by the low-level controller of our self-driving vehicle. We train our trajectory scoring model on a 500+ hour real-world dataset of expert driving demonstrations in Las Vegas within the maximum entropy IRL framework. DriveIRL's benefits include: a simple design due to only learning the trajectory scoring function, relatively interpretable features, and strong real-world performance. We validated DriveIRL on the Las Vegas Strip and demonstrated fully autonomous driving in heavy traffic, including scenarios involving cut-ins, abrupt braking by the lead vehicle, and hotel pickup/dropoff zones. Our dataset will be made public to help further research in this area.

## 1  Introduction

Self-driving cars have been the focus of significant research and development over the past decade. Some companies are tantalizingly close to deploying commercial self-driving taxi services that would make urban transportation cheaper and safer. Progress in self-driving cars has been largely driven by new datasets [Caesar et al., 2019, Sun et al., 2020, Chang et al., 2019, Geiger et al., 2012] that helped fuel dramatic improvements in machine learning approaches to object detection [Zhou and Tuzel, 2018, Lang et al., 2019] and motion forecasting [Cui et al., 2019, Chai et al., 2019, Phan-Minh et al., 2020]. However, the critical motion planning and decision-making algorithms that ultimately determine driving behavior have yet to see similar benefits from machine learning.

Classical planning and decision-making algorithms for self-driving cars rely heavily on hand-engineered components [Paden et al., 2016]. Developers will typically hand-tune the scoring function that determines which behaviors are desirable. Manually adjusting features and weights can be a painstaking process – improving performance in one area often causes unintended regressions elsewhere. Our planner avoids the need to manually craft detailed features or tune weights by learning these components from expert demonstrations using a maximum entropy IRL framework.
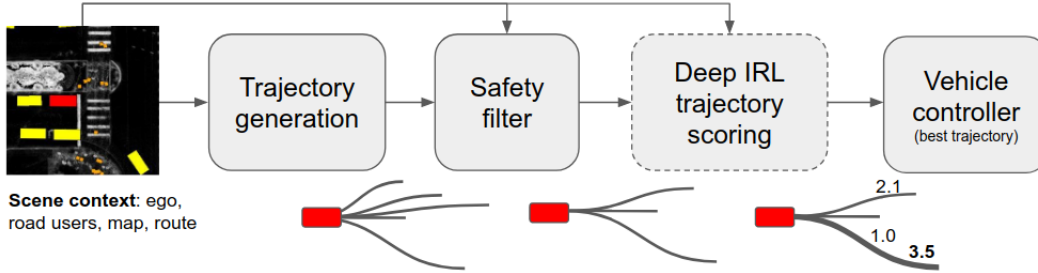
Figure 1: DriveIRL architecture. The learned scoring component is indicated with a dotted boundary.

Our DriveIRL system works by *generating*, *checking*, and *scoring* trajectories for our vehicle. We use simple and interpretable modules to do the relatively easy tasks of generating a diverse set of ego trajectories and checking that they are safe. Careful construction of the proposed trajectories ensures that they a) are dynamically feasible, b) follow the route, c) satisfy assumptions from the vehicle controller, and d) are diverse. We then apply a lightweight *safety filter* that ensures that each trajectory satisfies a recursive safety guarantee: if we execute the first part of the trajectory, there exists a safe continuation of that trajectory which avoids collision. The learning component of our model focuses entirely on appropriately scoring these trajectories based on the expert demonstrations. Our design directs the model capacity towards hard-to-specify nuances in behavior (e.g., speed profiles, clearances) instead of also creating "nice" trajectories and avoiding obviously unsafe behavior.

DriveIRL achieves strong real-world driving performance on the Las Vegas Strip. The Strip is a major thoroughfare in Las Vegas which connects many of the major hotels and casinos. Challenges include dense traffic, aggressive cut-ins, erratic drivers, and busy passenger pickup/dropoff zones near the hotels. We deployed DriveIRL on a self-driving car and drove fully autonomously on the Strip in these scenarios, showing the practical utility of our approach.

Our main contributions towards learning-based planning for self-driving cars are:

- The first learning-based planner to drive a car in dense, urban traffic using IRL.
- A simple yet powerful modeling framework that focuses learning on the aspect of driving that is most challenging to specify.
- Detailed evaluation of our planner on a real-world dataset that we will make public.

## 2 Related work

**Classical planning**: Traditional approaches formulate the planning problem as search over an appropriately constructed graph [e.g., A*, RRT*, PRM*; LaValle, 2006, Paden et al., 2016] or trajectory optimization [Paden et al., 2016]. These methods often have strong theoretical guarantees on convergence to an optimal solution and are relatively easy to interpret. However, the cost function that defines desired behavior is often hand-engineered, and in practice requires painstaking tuning to produce appropriate behavior.

**Imitation learning (IL)**: IL methods attempt to directly imitate the actions of an expert driver. It has seen applications to self-driving cars starting with the pioneering work of ALVINN [Pomerleau, 1988]. More recently, an end-to-end driving policy was learned from camera images to control actions for lane keeping [Bojarski et al., 2016].

A fundamental issue with IL is that there is a distribution shift from training to deployment, as small errors lead to the model operating outside of its training data, which then leads to larger errors. ChauffeurNet [Mayank Bansal, 2019] uses behavioral cloning with extensive data augmentation to mitigate the distribution shift issue, and UrbanDriver [Scheel et al., 2021] uses an offline policy gradient method with closed-loop rollouts during training to automatically create appropriate data augmentation. TrafficSim [Suo et al., 2021] similarly uses closed-loop training, but with a focus on creating traffic simulation. While data augmentation improves our performance, it is not as critical since our trajectory generation mechanism pulls the car towards the lane center, reducing divergence.

Closely related is work by Zeng et al. [2019] which learns a costmap over the environment to score a set of procedurally generated trajectories. Our approach improves on trajectory generation since we

ensure map compliance, as well as on scoring flexibility since we do not impose the assumption of an additive costmap. Furthermore, we demonstrate our model on a vehicle in dense urban traffic.

Another similar approach is that of Vitelli et al. [2021], where a hybrid model with a learned planner and an interpretable fallback layer drive in San Francisco. Our IRL-based model is simpler and less reliant on a fallback layer. Furthermore, the recursive check of our safety filter is less conservative.

**Reinforcement learning (RL)**: RL approaches learn a driving policy by optimizing a reward function. The standard approach requires a simulator [e.g., CARLA; Dosovitskiy et al., 2017] to update the environment that the driving policy interacts with. There have been a variety of approaches that have shown strong performance in simulation [Chen et al., 2020, 2021].

Real-world applications of RL for self-driving cars have been rarer, likely due to the difficulty in modeling the environment and specifying the reward function. An early notable example is Riedmiller et al. [2007], where they learn a steering policy for a real car. More recently, lane following was demonstrated using deep RL [Kendall et al., 2019]. This approach controlled both speed and steering on a real car. We contrast the rural driving evaluations above with our experiments in busy Las Vegas.

**Inverse reinforcement learning (IRL)**: IRL methods assume that the expert is optimizing an unknown cost function, which is learned from expert demonstrations. An early application of IRL to self-driving cars was for parking lot navigation [Abbeel and Ng, 2004]. The method learned multiple different driving styles from a handful of demonstrations. However, the environment was static and the formulation assumes a linear combination of carefully handcrafted features.

Our approach is based on the popular maximum entropy formulation of IRL [Ziebart et al., 2008], which avoids ambiguities inherent in matching feature expectations. The maximum entropy IRL approach was extended to deep learning in Wulfmeier et al. [2015], which avoided the need for laborious hand-engineering of features, and applied to simple benchmarks. The work of Huang et al. [2021] is the most related to our approach, but their model only learns a handful of feature weights while still assuming a linear combination of handcrafted features. In addition, their method is validated on a highway driving dataset and not on a real vehicle.

# 3 Inverse Reinforcement Learning Planner

In this section, we describe our Inverse Reinforcement Learning (IRL) Planner as shown in Fig. 1. Our system consists of three main stages: trajectory generation (Sec. 3.2), safety filtering (Sec. 3.3), and trajectory scoring (Sec. 3.4). We rely on simple and reliable hand-engineered modules for trajectory generation and safety, and focus on learning how to score trajectories.

## 3.1 Input and output

**Input**: We encode the environment (or scene) around our self-driving car using a mid-level representation. We assume that the ego is localized within a high-definition map and that objects are detected and tracked by a Perception system. Other road users (e.g., cars, bicyclists, and pedestrians) are represented by object type, an oriented bounding box, and speed. The high-definition map provides lane center-lines, road boundaries, traffic light locations, pedestrian crosswalks, speed limits, and other semantic information. We also provide a route, which indicates the lanes that the ego should traverse to make progress towards its goal.

We refer to the *scene context* at a given timestamp as a) the ego dynamic state $\mathcal{S}$ (speed, acceleration, steering), b) the other road users $\mathcal{U}$ (type, oriented bounding box, speed) , c) the map $\mathcal{M}$, and d) the ego's desired route $\mathcal{R}$. The model receives the *scene context* at the current timestamp as well as a specified number of previous timestamps (e.g., the past 1 second) as history $\mathcal{H}$.

**Output**: Our planner generates multiple ego trajectories and scores each one according to how closely it matches what an expert would do given the scene context. A trajectory is a discrete sequence of future states of the ego, where we assume that there is a fixed timestep between all states. Let $s_t = (x, y, \theta, v)$ represent a state at time $t$, with position $(x, y)$, heading $\theta$, and speed $v$. All values are with respect to the ego's geometric center in a fixed coordinate frame. The trajectory $\tau = [s_1, \ldots, s_T]$, where $T$ is the planned time horizon, that is ranked the best among a set of trajectories $\mathcal{T}$, is used as a reference for the vehicle's tracking and actuator controller.
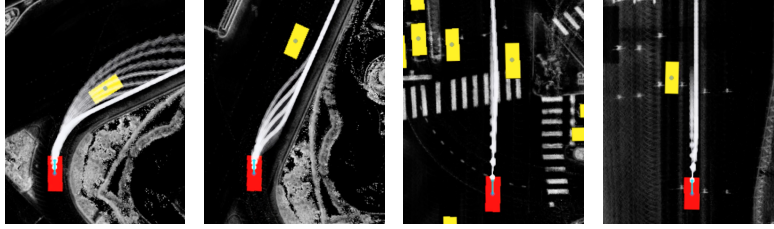
Figure 2: Proposed trajectories for the ego (red rectangle). Each trajectory is shown in translucent white dot-line. Overlap is due to multiple acceleration profiles. All trajectories return to the route.

## 3.2 Trajectory generation

The trajectory generation module uses the scene context to synthesize a diverse set of possible future motions for the ego. Important considerations for the ego's trajectory are that it a) is dynamically feasible, and b) satisfies all requirements of the low-level tracking and actuator control (i.e., levels of continuity, minimum turn radius, minimum acceleration from a stop). Secondary considerations are that the trajectory is compliant with the map (e.g., it stays on the road). While these considerations do not preclude using a learned trajectory generation module, we found that a hand-engineered trajectory generator best satisfied the considerations above.

The trajectory generator uses i) the current ego state $S$, ii) the route $R$, and iii) the map $M$ to create a diverse set of ego trajectories $T$, namely $(S, R, M) \mapsto T$. The generator integrates a desired acceleration profile along the route ahead of the ego. In our experiments, we specified a range of constant acceleration profiles ranging from a hard brake $(-5.0 \, \text{m/s}^2)$ to a moderate acceleration $(1.5 \, \text{m/s}^2)$. As the ego will not always be on the lane center-line (due to vehicle controller tracking errors), we smoothly connect the initial ego pose with the route with Dubins paths LaValle [2006] where turning radii are a fixed set of parameters. In a typical scene, the trajectory generator usually creates 50-150 trajectories depending on the ego state and route. Some examples are shown in the Fig 2. Appendix A.1 provides results to validate that the generated trajectories are of good quality.
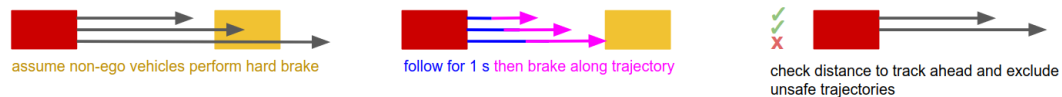
## 3.3 Safety filter



Figure 3: Safety filter. *Left*: toy scenario with three trajectories (ego is in red, vehicle ahead is in yellow). *Middle*: modified trajectories. *Right*: unsafe trajectories excluded from trajectory set.

Before scoring candidate trajectories, we apply an interpretable safety filter (Fig 3) to guarantee basic safety (i.e., no collisions). It consists of:

- a set of world assumptions used to predict the behavior of the non-ego road users,
- a set of trajectory modifiers which are applied to the ego trajectory, and
- a set of safety checks which the modified ego trajectory needs to pass.

For a candidate trajectory to be considered safe, it must pass all safety checks, under the given trajectory modifications and assumptions about the other road users. See A.2 for details.

Our safety filter is similar in spirit to the fallback layer proposed by Vitelli et al. [2021], except that 1) it directly filters the proposed trajectories, rather than projecting the output trajectory to an ad-hoc trajectory set, and 2) the trajectory modifier effectively implements a recursive safety guarantee with minimal assumptions and checks, without compromising comfort.

## 3.4 Trajectory scoring with maximum entropy IRL

Appropriately scoring trajectories is the core challenge of our planning approach. This difficulty is because proper driving behavior is heavily influenced by the environment around us, including other road user behavior and goals, of which we only have a partial understanding.

Trajectories are scored by a deep neural network trained with a maximum entropy IRL loss [Ziebart et al., 2008]. We use expert demonstrations collected from a skilled human driving our vehicle. The loss favors trajectories that most closely match the expert demonstration $\tau^\star$ in feature space. In particular, let $r(\tau)$ represent the reward of the trajectory $\tau \in \mathcal{T}$, the probability of a trajectory $\tau^*$ being selected according to the maximum entropy principle is $P(\tau^\star) = \frac{\exp r(\tau^\star)}{\sum_\tau \exp r(\tau)}$.

The negative log-likelihood loss (NLL) on a dataset $D$ is defined as $\ell(D) = -\sum_{d \in D} \log P(\tau^\star(d))$ where $\tau^\star(d)$ is the demonstrated trajectory on the token $d \in D$. To address data imbalance issues, we augment NLL with focal loss [Lin et al., 2017] (with a $\gamma$ of 2.0)

$$\ell(D) = -\sum_{d \in D} (1 - P(\tau^\star(d)))^\gamma \log P(\tau^\star(d)). \tag{1}$$

**Features**: We compute features for each proposed trajectory to use as inputs to our neural network. These features can be based on any combination of a proposed trajectory $\tau$, ego state $\mathcal{S}$, other road users $\mathcal{U}$, the map $\mathcal{M}$, route $\mathcal{R}$, and history $\mathcal{H}$, meaning that $F_i : (\tau, \mathcal{S}, \mathcal{U}, \mathcal{M}, \mathcal{R}, \mathcal{H}) \mapsto f_i \in \mathbb{R}^{k_i}$, where $F_i$ is the feature extraction function corresponding to feature $i$ and $k_i$ is its dimension.

- *Time-to-collision (TTC)*: the minimum number of seconds before the ego would collide with another road user in the (predicted) future. Evaluated at multiple points.
- *ACCInfo*: the ego speed, the distance to the road user ahead, the speed of the road user ahead, and the relative speed of the road user ahead. Evaluated at multiple points.
- *MaxJerk*: the maximum jerk $(\mathrm{m/s^3})$ along the trajectory.
- *MaxLateralAccel*: the max lateral acceleration $(\mathrm{m/s^2})$ along the trajectory.
- *PastCoupling*: concatenation of the future trajectory and the one second of past ego poses to model learn to maintain the coherence between the past, present, and future trajectories.
- *SpeedLimit*: how closely the trajectory obeys the speed limit. Evaluated at multiple points.

More implementation details can be found in the Appendix A.3.

**Motion prediction**: Some of the features computed for each proposed trajectory require an estimate of where other road users will be in the future, such as *Time-to-collision (TTC)* and *ACCInfo*. We use an Intelligent Driver Model (IDM) [Treiber et al., 2000] as our prediction model for other cars, with a conservative acceleration value to avoid assuming that stationary vehicles will speed up. We use a constant velocity model for pedestrians and for vehicles without a nearby lane.

**Model architecture**: To score a trajectory, we adopt an architecture in which the extracted features are processed separately before interacting with one another through a masked self-attention mechanism.
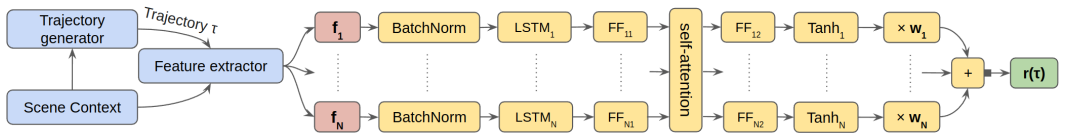


Figure 4: Detailed trajectory scoring architecture.

Under this architecture, each input feature $f_i$, as a temporal sequence of related vehicle-environment interaction data, is first normalized through an application of a BatchNorm1D layer before being fed to an LSTM module with one layer and a hidden size of 20. The output of the LSTM becomes the input to a feed-forward module and then a self-attention mechanism with two heads and an embedding dimension of 120. Here we employ zero-masking of the queries to encode position. By taking into account other features through self-attention, the model produces for each feature a "corrected" output embedding that can now be passed to a feed-forward network which converts it into a scalar and then a $\tanh$ activation to produce a feature score $y_i$. The final score for the trajectory is

Table 1: A detailed distribution of our dataset (182, 032 total 11-second scenarios).

| Tags | Straight | Right | Left | Stopped | Slow | Intersection | Close | ASV |
|------|----------|-------|------|---------|------|--------------|-------|-----|
| Scenarios | $163,079$ | $1,463$ | $2,578$ | $53,566$ | $37,414$ | $28,197$ | $11,232$ | $15,329$ |
| Ratio | $89.6\%$ | $0.8\%$ | $1.4\%$ | $29.4\%$ | $20.6\%$ | $15.5\%$ | $6.2\%$ | $8.4\%$ |

the sum of these feature scores after they are multiplied by the corresponding learnable feature weight parameters $w_i$: $r(\tau) = \sum_i w_i y_i$. In total, our base (best) model has $\approx 88,700$ trainable parameters.

## 4  Experiments

The proposed inverse reinforcement leaning planner was evaluated on a large-scale dataset and the results are presented in the following. The dataset we used is described in Sec 4.1. The metrics for comparison is explained in Sec 4.2. Various model ablation studies and the comparison with baseline are shown in Sec 4.3 and 4.4. We demonstrated both simulation results and the real-world driving tests in Sec 4.5 and 4.6 respectively.

### 4.1  Dataset

We created a self-driving car dataset that captures real-world urban driving in the center of Las Vegas. Our dataset is a part of the nuPlan [Caesar et al., 2021] dataset that will be made public. It includes object annotations and high-definition maps. Vehicles, pedestrians, and bicyclists are automatically annotated using an offline perception system (similar in spirit to Qi et al. [2021]) and viewed as ground truth. We performed filtering and extracted 182,032 scenarios, each 11 seconds in duration (1 second past, 10 seconds future), for a total of approximately 556 hours. Our main interest was to learn good adaptive cruise control (ACC) behavior. Thus, we filtered out scenarios where the ego made lane changes or deviated far from the lane. After filtering, we performed a 3:1:1 split for train, val, and test sets. Tab. 1 shows a detailed distribution of our dataset by scenario tags. The tags in the table are not mutually exclusive and a scenario can belong to multiple tags. More detailed definitions for the scenario tags are in Appendix A.4.

### 4.2  Metrics

We evaluate our model using a variety of metrics to give a full picture of driving. To approximate real-world conditions, we perform a closed-loop replay for each scenario for a duration of 10 seconds. We initialize the ego at the start of the scene, compute a planned trajectory, move along that trajectory for one step, replay the other agents, and repeat. Then, we compute metrics on the resulting executed trajectory as averages over the full scene duration.

**Metric computation**: Evaluation was done with a time step size of 0.2 seconds and a total duration of 10 seconds. The model was given 1 second of ground truth past prior to the start of the scene. Other road users were updated by replaying their positions from the database.

**Metric categories**: We have four high level categories of metrics that contain "low level" metrics and high level summary metrics that act as a score for the category. The categories for our metrics are Safety, Comfort, Progress, and $\ell_2$ (with a yaw penalty of 2.5). Further details about our metric categories and the low level computations that make them up can be found in Appendix A.5.

**Metric limitations**: Currently, there are two major limitations to our metrics evaluation.

*The use of replay for other agents.* Since other vehicles do not react to the ego (e.g., if we drive slower than the expert in the data), the overall "Safety" score is a lower bound on safety.

*No controller or vehicle dynamic simulation for the Ego.* We currently "teleport" the ego along its trajectory, causing jerk to be erroneously high in some cases. Similar to safety, this makes our "Comfort" score a lower bound on what we actually observe in real world driving.

Table 2: Ablation study on the importance of each feature. The row indicates the feature removed from the baseline model. See Sec. 3.4 for definitions.

| Model | Safety | Comfort | Progress | $\ell_2$ (w/ yaw) | Collision | Tailgate |
|---|---|---|---|---|---|---|
| Base (ours) | **0.925** | **0.840** | **0.988** | 2.290 | **0.001** | 0.015 |
| No TTC | 0.865 | 0.790 | 0.958 | 2.679 | 0.055 | 0.070 |
| No ACCInfo | 0.862 | 0.817 | 0.959 | 2.832 | 0.004 | 0.035 |
| No MaxJerk | 0.863 | 0.825 | 0.960 | **2.201** | **0.001** | 0.029 |
| No MaxLatAccel | 0.917 | 0.821 | 0.982 | 2.299 | 0.005 | **0.011** |
| No PastCoupling | 0.901 | 0.697 | 0.979 | 2.905 | 0.005 | **0.011** |
| No SpeedLimit | 0.881 | 0.809 | 0.987 | 2.483 | 0.002 | 0.028 |

Table 3: Comparison between different augmentation schemes.

| Model | Safety | Comfort | Progress | $\ell_2$ (w/ yaw) | Collision | Tailgate |
|---|---|---|---|---|---|---|
| Base (low noise) | **0.925** | 0.840 | **0.988** | 2.290 | **0.001** | **0.015** |
| No noise | 0.850 | **0.850** | 0.956 | **2.289** | 0.005 | 0.055 |
| High noise | 0.917 | 0.845 | 0.986 | 2.525 | 0.002 | 0.016 |
| Low past + present | 0.921 | 0.817 | 0.982 | 2.302 | 0.005 | 0.019 |

## 4.3 Model ablations

In our experiments, we use a batch size of $64$ and an Adam optimizer with an initial learning rate of $10^{-3}$. Additionally, we use a "cosine annealing with warm restarts" scheduler, which gradually lowers the learning rate to a minimum of $10^{-4}$ and resets it every seven epochs. All models are trained over 20 epochs on eight AWS g4dn-metal instances with eight 16 GB NVIDIA Tesla T4 GPUs each. Because closed-loop simulation is computationally expensive, we randomly sampled $1,000$ scenarios from our evaluation set for ablation studies and $3,000$ scenarios from the test set for the final performance evaluation against other baselines. Training and closed-loop metrics evaluation takes about an hour per epoch.

**Feature importance**: To understand the importance of each hand-engineered feature and the main contribution of each, an ablation study for features is conducted and summarized in Tab. 2. The relative importance of each feature is shown by dropping one of them out at a time. We claim that all the features are important because the *Base* model which includes all features got the highest scores across all high-level metrics and had lowest Collision rate. Even though the $\ell_2$ error is a bit higher compared to *No MaxJerk*, the $0.089\,m$ difference is not significant in the qualitative results. The results also demonstrated the importance of *PastCoupling* feature in ensuring Comfort. The experiment also showed that *TTC* feature contributes significantly to reducing the collision rate.

**Data augmentation**: Data augmentation is important to ensure that our model can learn how to recover from errors. Since the reference trajectory is never followed perfectly by the vehicle, errors can accumulate. We perturb the ego's initial state during training to reduce the sensitivity to such errors. For our low noise baseline, we use zero-mean Gaussian data augmentation for longitudinal offset ($1.2\,m$ std), lateral offset ($0.8\,m$ std), heading offset ($0.1\,rad$ std), and velocity ($0.1\,m/s$ std). For the high noise ablation, we respectively use $2.5\,m$ std, $1.5\,m$ std, $0.3\,rad$ std, and $0.2\,m/s$ std. We clamp velocity to avoid negative values. Several example images are shown in A.6.

**Architecture**: We perform several ablations on the model architecture before selecting an architecture in which the extracted features are processed separately before interacting with one another through a masked self-attention mechanism. We show in Tab. 4 that the other two extremes, namely, concatenating all input features and using them as one monolithic feature in a single feedforward network or siloing all input features (not allowing any interaction through attention or otherwise) have resulted in inferior performance. It is also seen that input normalization and attention input masking are beneficial.

**Loss**: Tab. 5 shows that it is better to maximize the probability the projection of the ground truth onto the trajectory set (the best approximation in average $\ell_2$ norm) instead of the ground truth itself. This makes sense because the ground truth does not come from the same distribution as the proposals and

Table 4: Comparison between different model architectures.

| Model | Safety | Comfort | Progress | $\ell_2$ (w/ yaw) | Collision | Tailgate |
|---|---|---|---|---|---|---|
| Base (ours) | **0.925** | **0.840** | **0.988** | 2.290 | **0.001** | 0.015 |
| Monofeature FC | 0.816 | 0.715 | 0.922 | 3.084 | 0.009 | 0.020 |
| Siloed FC | 0.900 | 0.784 | 0.964 | 2.394 | 0.018 | 0.021 |
| No input norm | 0.880 | 0.817 | 0.957 | 2.496 | 0.003 | 0.016 |
| Attention no masking | 0.910 | 0.835 | 0.983 | **2.285** | 0.004 | **0.012** |

Table 5: Comparison between different loss functions.

| Metric | Safety | Comfort | Progress | $\ell_2$ (w/ yaw) | Collision | Tailgate |
|---|---|---|---|---|---|---|
| Base (ours) | 0.925 | **0.840** | 0.988 | **2.290** | **0.001** | **0.015** |
| GT as demo | 0.910 | 0.553 | **0.992** | 2.870 | 0.012 | 0.029 |
| Possibly unsafe demo | 0.873 | 0.823 | 0.977 | 2.518 | 0.036 | 0.049 |
| Demo w/ weighted yaw | **0.932** | 0.839 | 0.984 | 2.530 | 0.005 | 0.018 |
| Without focal loss | 0.910 | 0.831 | 0.976 | 2.393 | 0.004 | 0.018 |

is not available at inference time. Filtering possibly unsafe trajectories from the set before finding the ground truth projection is also crucial to obtaining a safe model. Doing the projection using the average $\ell_2$ norm instead of an $\ell_2$ norm with a yaw error penalty also seems favorable. Lastly from the same table, we can see that using focal loss as in Equation 1 improves performance. Another experiment in Appendix A.7 that compares focal loss against training on a better balanced dataset also shows that using focal loss is actually more effective for DriveIRL.

## 4.4 Baselines

In this section, we evaluate our model on a test dataset and compare it with an Intelligent Driver Model (IDM) [Treiber et al., 2000] and a constant speed (CS) lane follow model. The IDM baseline is a reasonable choice because it is a well-known version of an expert planner that focuses on adaptive cruise control. Meanwhile, the CS lane follow model is a simple lower-bound. The results are shown in Tab. 6. Our base model plus safety filter outperforms others in all safety related metrics, and that shows the safety filter protects the vehicle on several collision cases our model cannot handle perfectly. Without the safety filter, our base model still outperforms the IDM baseline. The IDM model has significantly higher $\ell_2$ error, indicating that the IDM model does not drive like a human expert. Furthermore, our base model also has higher scores in all safety related metrics in both high- and low-level scores like collision rate and tailgate rate.

## 4.5 Simulation results

Fig. 5 exhibits some qualitative closed-loop simulation results of our planner driving in typical scenarios. These scenarios are shown as a sequence of snapshots along the closed-loop rollout. The ego vehicle is shown as a red rectangle, the expert vehicle is in blue, and other vehicles are in yellow. The orange line is the planned route (along the lane centerline) and the purple circles represent the

Table 6: Baselines on the test set. IDM = Intelligent Driver Model. CS = constant speed.

| Metric | Safety | Comfort | Progress | $\ell_2$ (w/ yaw) | Collision | Tailgate |
|---|---|---|---|---|---|---|
| Expert | 1.000 | 0.984 | 1.000 | 0.000 | 0.000 | 0.000 |
| Base + Safety (ours) | **0.930** | 0.815 | 0.971 | **2.204** | **0.001** | **0.006** |
| Base (ours) | 0.916 | 0.830 | 0.986 | 2.351 | 0.003 | 0.017 |
| IDM | 0.891 | 0.898 | **0.987** | 4.478 | 0.005 | 0.019 |
| CS lane follow | 0.669 | **0.992** | 0.902 | 3.963 | 0.161 | 0.166 |

(a) Ego starting from a stop.

(b) Ego stopping for the lead vehicle.



(c) Adaptive cruise control.
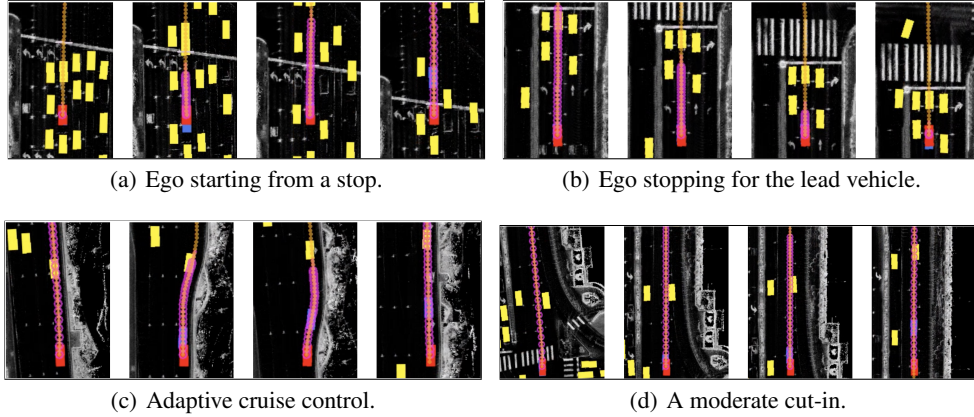
(d) A moderate cut-in.

Figure 5: Qualitative driving performance in common scenarios. Full-size images are in A.8.



Figure 6: Smoothly stopping behind a vehicle in dense traffic on the Las Vegas Strip.

planned trajectory for the next 6 seconds. Our planner shows good performance over a range of scenarios, exhibiting reasonable and consistent behavior over 10 second rollouts.

## 4.6 Real-world driving

Prior to deploying on public roads, DriveIRL was rigorously tested in both simulation and on private, closed-course routes. The simulation tests consist of the same Las Vegas Strip route that was our deployment goal, and involve a high-fidelity dynamics model for the ego vehicle and numerous actors exhibiting a wide variety of behaviors. When deployed on the Strip, the vehicle was piloted by a vehicle operator who was trained to take over for unsafe behavior and situations outside of our operating domain, including construction zones, bus stops, and yielding for emergency vehicles.

On the Strip, our planner handled challenging scenarios such as heavy traffic, aggressive cut-ins, unpredictable drivers, and busy passenger pick-up/drop-off zones near the hotels and casinos. Without the safety filter, the vehicle remained in autonomous mode for 8.8 miles of the 11-mile route. Overrides occurred for mandatory takeover regions and twice for undesired behavior. With the safety filter, the vehicle remained in autonomous mode for 6.9 of 8.5 miles, with takeovers only occurring due to mandatory takeover regions.

Fig. 6 shows a typical maneuver where our self-driving vehicle smoothly stops for a vehicle ahead while surrounded by multiple vehicles. In Sec. A.9, we have included video clips with more real-world driving maneuvers. These videos are grouped in categories such as cut-ins, driving around passenger pickup/dropoff zones, driving with a vehicle ahead and stopping behind a vehicle.

## 5 Conclusions

We introduced DriveIRL: the first learning-based planner to control a car in dense, urban traffic using inverse reinforcement learning. By designing an architecture split into ego trajectory *generation*, *checking*, and *scoring*, we were able to leverage simple and reliable methods for the relatively easy tasks of trajectory generation and safety checking. This architecture allowed the trajectory scoring component of our system to focus on learning nuanced driving behavior important for good

performance in dense traffic. We demonstrated our planner on the busy Las Vegas Strip, where it showed strong real-world performance on challenging scenarios such as cut-ins, abrupt braking, and cluttered hotel pickup/dropoff zones.

## Acknowledgments and Disclosure of Funding

## References

Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuScenes: A multimodal dataset for autonomous driving. *arXiv preprint arXiv:1903.11027*, 2019.

Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in perception for autonomous driving: Waymo Open Dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, and James Hays. Argoverse: 3d tracking and forecasting with rich maps. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012.

Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4490–4499, 2018.

Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12697–12705, 2019.

H. Cui, V. Radosavljevic, F. Chou, T. Lin, T. Nguyen, T. Huang, J. Schneider, and N. Djuric. Multimodal trajectory predictions for autonomous driving using deep convolutional networks. In *International Conference on Robotics and Automation (ICRA)*, May 2019.

Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. MultiPath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. In *3rd Conference on Robot Learning (CoRL)*, November 2019.

Tung Phan-Minh, Elena Corina Grigore, Freddy A Boulton, Oscar Beijbom, and Eric M Wolff. Covernet: Multimodal behavior prediction using trajectory sets. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14074–14083, 2020.

Brian Paden, Michal Čáp, Sze Zheng Yong, Dmitry Yershov, and Emilio Frazzoli. A survey of motion planning and control techniques for self-driving urban vehicles. *IEEE Transactions on intelligent vehicles*, 1(1):33–55, 2016.

Steven M. LaValle. *Planning Algorithms*. Cambridge University Press, New York, NY, USA, 2006. ISBN 0521862051.

Dean A Pomerleau. Alvinn: An autonomous land vehicle in a neural network. *Advances in neural information processing systems*, 1, 1988.

Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Prasoon Goyal, Lawrence D. Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, Xin Zhang, Jake Zhao, and Karol Zieba. End to end learning for self-driving cars, 2016. URL `https://arxiv.org/abs/1604.07316`.

Abhijit Ogale Mayank Bansal, Alex Krizhevsky. ChauffeurNet: Learning to drive by imitating the best and synthesizing the worst. In *Robotics: Science and Systems (RSS)*, June 2019.

Oliver Scheel, Luca Bergamini, Maciej Wołczyk, Błażej Osiński, and Peter Ondruska. Urban driver: Learning to drive from real-world demonstrations using policy gradients, 2021. URL `https://arxiv.org/abs/2109.13333`.

Simon Suo, Sebastian Regalado, Sergio Casas, and Raquel Urtasun. TrafficSim: Learning to simulate realistic multi-agent behaviors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10400–10409, 2021.

Wenyuan Zeng, Wenjie Luo, Simon Suo, Abbas Sadat, Bin Yang, Sergio Casas, and Raquel Urtasun. End-to-end interpretable neural motion planner. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

Matt Vitelli, Yan Chang, Yawei Ye, Maciej Wołczyk, Błażej Osiński, Moritz Niendorf, Hugo Grimmett, Qiangui Huang, Ashesh Jain, and Peter Ondruska. Safetynet: Safe planning for real-world self-driving vehicles using machine-learned policies, 2021. URL `https://arxiv.org/abs/2109.13602`.

Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. In *Conference on robot learning*, pages 1–16. PMLR, 2017.

Dian Chen, Brady Zhou, Vladlen Koltun, and Philipp Krähenbühl. Learning by cheating. In *Conference on Robot Learning*, pages 66–75. PMLR, 2020.

Dian Chen, Vladlen Koltun, and Philipp Krähenbühl. Learning to drive from a world on rails. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15590–15599, 2021.

Martin Riedmiller, Mike Montemerlo, and Hendrik Dahlkamp. Learning to drive a real car in 20 minutes. In *2007 Frontiers in the Convergence of Bioscience and Information Technologies*, pages 645–650. IEEE, 2007.

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 *2019 International Conference on Robotics and Automation (ICRA)*, pages 8248–8254. IEEE, 2019.

Pieter Abbeel and Andrew Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *Proceedings of the Twenty-First International Conference on Machine Learning*, ICML '04, page 1, New York, NY, USA, 2004. Association for Computing Machinery. URL `https://doi.org/10.1145/1015330.1015430`.

Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse reinforcement learning. In *Aaai*, volume 8, pages 1433–1438. Chicago, IL, USA, 2008.

Markus Wulfmeier, Peter Ondruska, and Ingmar Posner. Maximum entropy deep inverse reinforcement learning. *arXiv preprint arXiv:1507.04888*, 2015.

Zhiyu Huang, Jingda Wu, and Chen Lv. Driving behavior modeling using naturalistic human driving data with inverse reinforcement learning. *IEEE Transactions on Intelligent Transportation Systems*, pages 1–13, 2021. doi: 10.1109/TITS.2021.3088935.

Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

Martin Treiber, Ansgar Hennecke, and Dirk Helbing. Congested traffic states in empirical observations and microscopic simulations. *Physical review E*, 62(2):1805, 2000.

Holger Caesar, Juraj Kabzan, Kok Seang Tan, Whye Kit Fong, Eric Wolff, Alex Lang, Luke Fletcher, Oscar Beijbom, and Sammy Omari. nuplan: A closed-loop ml-based planning benchmark for autonomous vehicles. *arXiv preprint arXiv:2106.11810*, 2021.

Charles R. Qi, Yin Zhou, Mahyar Najibi, Pei Sun, Khoa Vo, Boyang Deng, and Dragomir Anguelov. Offboard 3d object detection from point cloud sequences. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6134–6144, June 2021.

# A  Supplemental Material

## A.1  Trajectory quality

We evaluate the quality of the trajectory generation by projecting the expert (ego) trajectory on the closest (in $\ell_2$) trajectory from the trajectory generator (Expert Projection). We compare this with the actual expert trajectory (Expert) and the trajectory chosen by our best model (Base) on 6-second scenarios in Tab. 7. Note the use of shorter scenarios to match the length of the (open-loop) trajectory. The low $\ell_2$ for the Expert Projection indicates that there is enough diversity in the trajectory set to closely match the expert. These results also indicate that, given the trajectory set, our model is near ceiling performance in terms of the high-level metrics, but it could improve on $\ell_2$.

Table 7: Trajectory set quality.

| Metric | Safety | Comfort | Progress | $\ell_2$ (w/ yaw) | Collision | Tailgate |
|---|---|---|---|---|---|---|
| Base (ours) | 0.968 | 0.851 | 0.989 | 1.216 | 0.000 | 0.004 |
| Expert | 1.000 | 0.962 | 1.000 | 0.000 | 0.000 | 0.000 |
| Expert Projection | 0.984 | 0.880 | 0.996 | 0.309 | 0.000 | 0.000 |

## A.2  Safety filter details

For each candidate trajectory, we check that the distance to the track ahead never falls below $1.5\,\mathrm{m}$, assuming that all non-ego vehicles perform a hard brake at $3.5\,\mathrm{m/s^2}$. We perform the check on modified ego trajectories, where each trajectory is followed for the first $1\,\mathrm{s}$, after which firm braking is applied along the trajectory at $2.5\,\mathrm{m/s^2}$, subject to a $3.5\,\mathrm{m/s^3}$ jerk limitation. In effect, this implements a kind of recursive safety guarantee: if we follow the trajectory for $1\,\mathrm{s}$ and then brake, are we still safe?

## A.3  Feature implementation details

We compute features for each proposed trajectory to use as inputs to our neural network. These features are represented as $F^i(\tau, \mathcal{S}, \mathcal{U}, \mathcal{M}, \mathcal{R}, \mathcal{H}) = f^i \in \mathbb{R}^{k_i}$, where $F^i$ is the feature extraction function corresponding to feature $i$ and $k_i$ is its dimension, and the output feature vector is $f^i$. A feature vector $f^i$ can be either a sequence of feature tuples $[f_1, \ldots, f_T]$ where $f_t$ is a feature tuple at timestamp $t$, or just a single feature tuple $[f_w]$ which encodes the characteristics of the whole proposed trajectory regardless of timestamps.

- *Time-to-collision (TTC)*: The number of seconds before (predicted) collision with other agents. Because the computation time for this feature is expensive, we sub-sample the trajectory waypoints at $[0.2, 0.4, 0.6, 1.0, 2.0, 4.0]$ seconds. The TTC saturation time is set to $4.0$ seconds, meaning agents that the ego is not projected to collide into within $4.0$ seconds are effectively ignored. Namely, the feature function is: $F_{TTC}(\tau, \mathcal{S}, \mathcal{U}, \mathcal{M}) \to [f_{1\_collide}, \ldots, f_{T\_collide}]$ where $f_{t\_collide}$ denotes the minimum time to collision at timestamp $t$.

- *ACCInfo*: The information needed for adaptive cruise control. This feature complements the time-to-collision feature. When ego is slowly moving forward, the time-to-collision might not be very meaningful as its value may still be large even when the ego is already very close to the vehicle in front, which is not desirable. *ACCInfo* directly addresses this issue by feeding the distance information to the model along with some other information. The feature function is $F_{ACCInfo}(\tau, \mathcal{S}, \mathcal{U}, \mathcal{M}) \to [f_{1\_acc}, \ldots, f_{T\_acc}]$ where $f_{t\_acc}$ is $(d_{ahead}, b_{ahead}, v_{ego}, v_{ahead}, v_{ego} - v_{ahead})_t$ where $d_{ahead}$ is the distance to the track ahead, $b_{ahead}$ is a Boolean flag indicating if there is a track in front given a tunable distance, which by default is $20\,\mathrm{m}$, $v_{ego}$ and $v_{ahead}$ are the speed of the ego and the track ahead respectively, and lastly the relative speed between the ego and the track ahead is added.

- *MaxJerk*: The maximum jerk $(\mathrm{m/s^3})$ along the trajectory. This feature does not just compute the jerk value at the current time given the past ego states. Rather, it computes all jerk values along the proposed trajectory, and then takes a maximum among these

values as the single max jerk value $j_{\max}$. This feature function can be represented as: $F_{MaxJerk}(\tau, \mathcal{H}) \to [f_{w\_jerk}, j_{\max}]$. The way that we represent the jerk value in this feature tuple $f_{w\_jerk}$ is through a series of Boolean flags as this representation is easier for the model to consume, similar to the idea of one-hot encoding. Given a range of interest for the jerk value, a set of thresholds are assigned among the range with a given resolution. The range in our experiments used is $[0, 10]$ with a fixed step size of $0.5\,\mathrm{m/s^3}$. If the jerk value is smaller than a threshold, the corresponding flag will be set to $1$; otherwise $0$, meaning smaller max jerk value gets more flags set to $1$. We also append the original jerk value $j_{\max}$ at the end of this feature tuple for the model to consume.

- *MaxLateralAccel*: The max lateral acceleration $(\mathrm{m/s^2})$ along the trajectory. Similar to the MaxJerk feature, the MaxLateralAccel feature computes the maximum acceleration in a lateral direction w.r.t. the ego, representing another measurement of comfort. The same quantization representation as used for MaxJerk above is used in this feature as well, namely $F_{MaxLateralAccel}(\tau, \mathcal{H}) \to [f_{w\_lat}, a^{lat}_{\max}]$ where $f_{w\_lat}$ is a series of Boolean flags. The threshold values range from $0$ to $5\,\mathrm{m/s^2}$ in step sizes of $0.2\,\mathrm{m/s^2}$. Note that original value of the maximum lateral acceleration $a^{lat}_{\max}$ is appended at the end as well.

- *PastCoupling*: The concatenation of the proposed trajectory and the past ego states. As another complementary comfort feature, we found that a feature that represents all information across all time stamps is also helpful in terms of comfort level because it makes model learn the level of coherence between the past, present, and future trajectory. Specifically, $F_{PastCoupling}(\tau, \mathcal{H}) \to [f_{1\_coup}, \ldots, f_{T\_coup}]$ where $f_{t\_coup}$ is a five-tuple $(x, y, \theta, v, a)_t$ at the timestamp $t$, and $v$ and $a$ stand for the ego speed and acceleration.

- *SpeedLimit*: Whether the ego drives under or close to the speed limit. The speed limit of the road block can be queried from our map data given the ego position. For each waypoint in our proposed trajectory, two values are computed, a relative speed based on the difference between the speed limit and the current speed, and a Boolean flag $b_{limit}$ to indicate if the ego speed is over the speed limit, meaning that the feature function is: $F_{SpeedLimit}(\tau, \mathcal{S}, \mathcal{M}) \to [f_{1\_limit}, \ldots, f_{T\_limit}]$ where $f_{t\_limit}$ is a two tuple $((v_{ego} - v_{limit})/v_{limit}, b_{limit})_t$. To make the relative speed more uniform between the high speed case and the low speed case, the relative speed is also normalized by the value of speed limit itself.

## A.4 Scenario tags

A detailed distribution of our dataset is shown in Tab. 1. The scenario tags in the table are not mutually exclusive and a sample can belong to multiple tags.

We classified a scenario as *Straight* if the ego drove straight throughout the scenario and the changes in the initial and final heading angles of ego were less than $0.1\,\mathrm{rad}$. A scenario was classified as a *Right* or *Left* turn if the ego made a turn in the scenario and the change in the heading angles was greater than $\frac{\pi}{3}\,\mathrm{rad}$. Due to the geometry of the Las Vegas Strip, the ego mostly drove straight. Turns were made occasionally to enter or leave pick-up or drop-off areas. As the definitions in *Straight* and *Turn* tags indicate, these tags do not form a perfect partition. Thus, the number of scenarios for these tags do not add up to the total.

*Stopped* means the ego did not move through out the scenario. This is largely because the ego was in the middle of a traffic jam or waiting for pedestrians. We classified a scenario as *Slow* if the ego was driving slower than $2.64\,\mathrm{m/s}$. In our dataset, ego drove relatively slowly. For example, we had $33,646$ scenarios where ego drove faster than $9.64\,\mathrm{m/s}$. This is because of the following reasons: i) the speed limit of the Strip was $15\,\mathrm{m/s}$ and ii) there was a lot of traffic in the Strip. The heavy traffic makes our dataset low-speed but challenging.

Whenever the ego went through an intersection area during a scenario, it was tagged as *Intersections* regardless of the traffic light status.

There are also tags that capture the interactions of the ego and the leading vehicle, namely *Close* and *Approaching Stopped Vehicle (ASV)*. A scenario was classified as *Close* if the ego had time gap less than $1.7\,\mathrm{s}$, calculated by distance to the vehicle ahead divided by ego speed. *ASV* was when the ego was moving toward a leading stopped vehicle and got closer than $10\,\mathrm{m}$.
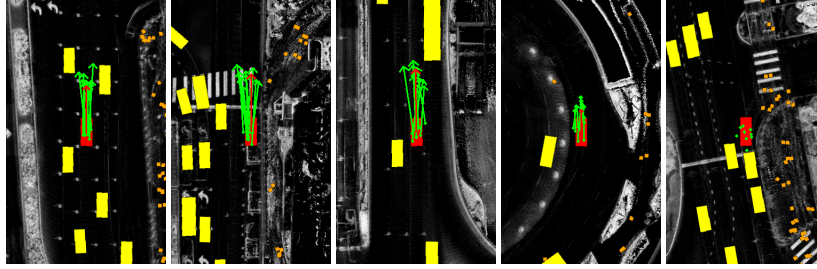
Figure 7: Data augmentation examples.

## A.5 Metrics breakdown

The following is a break down of our metric categories and how their respective metrics are computed:

- *Safety*: is the most important category and it is composed of the following items:
  - *collision rate*: Calculated as whether the ego hit any vehicle in *front* of it. [1]
  - *ego not off-road*: Measures whether or not the ego vehicle goes off the road at any point of the scene.
  - *ego minimum time to collision* $> 0.95\,\mathrm{s}$ [2]: Measures whether or not the time to collision of the ego was always above the set threshold
  - *tailgate rate, i.e. ego distance to lead vehicle* $< 1.5\,\mathrm{m}$ : Measures whether the ego was appropriately far from lead vehicles. Note that this metric is especially important during slow scenarios where the time to collision is not very informative.

- *Comfort*: is based on acceleration, jerk and yaw rate and is considering passing based on whether they are within the following thresholds: $-4.05\,\mathrm{m/s^2} <$ *ego longitudinal acceleration* $< 2.40\,\mathrm{m/s^2}$, *ego lateral acceleration* $< 4.89\,\mathrm{m/s^2}$, *ego absolute yaw rate* $< 0.95\,\mathrm{rad/s}$, *ego absolute yaw acceleration* $< 1.93\,\mathrm{rad/s^2}$, *ego absolute longitudinal jerk* $< 4.13\,\mathrm{m/s^3}$, and *ego absolute jerk magnitude* $< 8.37\,\mathrm{m/s^3}$.

- *Progress*: is made up of only two components, (1) whether or not the ego has made positive distance to the goal of more then $1\,\mathrm{m}$ when projected onto the route given to the ego and (2) whether or not the ego has diverged from the route, measured as whether the maximum distance from the ego to the route is greater than $4\,\mathrm{m}$.

- $\ell_2$: These metrics measure the distance between the resulting trajectory of the ego and the path the expert actually drove on the scenario. They have no high level summary. Some examples of $\ell_2$ metrics that have been useful for us include *ego average l2 with yaw* and *ego average cross-track (i.e. longitudinal) l2*. For $\ell_2$ with yaw this is calculated as standard $\ell_2$ between the trajectories with an adjustment of $2.5$ times the principal value of the heading differences between the trajectories.

## A.6 Data augmentation examples

Several examples of data augmentation are shown in Fig. 7. The original ego pose and velocity are shown as a red rectangle and arrow respectively while the speed scales with the length of arrow. The augmented ego center positions and velocities are shown as green points and arrows respectively. We plot 10 random results for demonstration.

---

[1]We use front collision rather than total collision, because we are replaying the other agents, which leads to many false positive rear collisions that are not caused by the ego's behavior. Specifically, if the ego stops but the expert did not then cars that were following the expert may then proceed to collide with the ego vehicle.

[2]All thresholds for metrics were derived from examining expert driving and were adjusted based on feedback of what worked and what didn't during deployment in the Las Vegas Strip.

Table 8: A detailed distribution of our balanced train set ($125, 875$ total scenarios).

| Tags | Straight | Right | Left | Stopped | Slow | Intersection | Close | ASV |
|---|---|---|---|---|---|---|---|---|
| Scenarios | $103,411$ | $3,340$ | $4,464$ | $32,273$ | $26,339$ | $24,784$ | $12,600$ | $12,287$ |
| Ratio | $82.2\%$ | $2.7\%$ | $3.5\%$ | $25.6\%$ | $20.9\%$ | $19.7\%$ | $10.0\%$ | $9.8\%$ |

Table 9: Data balancing.

| Metric | Safety | Comfort | Progress | $\ell_2$ (with yaw) | Collision rate | Tailgate rate |
|---|---|---|---|---|---|---|
| Base (ours) | **0.925** | **0.840** | **0.988** | 2.290 | **0.001** | **0.015** |
| Balanced set + base | 0.913 | 0.828 | 0.984 | **2.209** | 0.005 | 0.022 |
| Without focal | 0.910 | 0.831 | 0.976 | 2.393 | 0.004 | 0.018 |
| Balanced set + without focal | **0.922** | **0.851** | **0.981** | **2.280** | **0.002** | **0.014** |

## A.7 Data balancing

As shown in Tab. 1, the distribution of our dataset is unbalanced. For example, there are a lot more *Straight* scenarios compared to *Right* or *Left* turn scenarios. In addition, the interaction between the ego and the leading vehicle is crucial in ACC. This interaction is captured in *Close* and *ASV* scenarios, but they account for only $6.2\%$ and $8.4\%$, respectively. We performed balancing on our train set to make sure our model learns rare but important scenarios well. To be more specific, we took an upsampling approach and increased the ratio of rare scenarios by duplicating them multiple times in the train set. Even though the scenarios were copied, our data augmentation gave these scenarios some variation. Tab. 8 shows the distribution of our balanced train set. Compared to Tab. 1, the ratios of rare scenarios such as *Right*, *Left*, *Close* and *ASV* increased whereas those of frequent scenarios such as *Stopped* and *Straight* decrease. We performed balancing only on the train set and did not modify the validation and test sets.

The first two rows of Tab. 9 enables a comparison between using the original (unbalanced) train set and the balanced one. The next two rows are for comparing two train sets when trained without local loss. The balanced train set did not help much when applied on our base model. However, the balanced train set showed improvement when applied on the model without focal loss. This indicates using focal loss has a similar effect as using a balanced train set as expected. In fact, by comparing the first (with focal loss + unbalanced train set) and fourth (without focal loss + balanced train set) rows, we see that using focal loss and a balanced train set have comparable performances. Tab. 10 compares the $\ell_2$ (with yaw) metrics of each scenario tags for the three models of *Without focal loss*, *Base (ours)*, and *Balanced train set + without focal loss*. We can see that the base model that used focal loss achieved better $\ell_2$ (with yaw) for rare scenarios such as *Right*, *Left*, and *Close*. This suggests using the focal loss is more effective in resolving the dataset imbalance issue.

## A.8 Simulation video clips

Tab. 11 presents a curated list of video clips where we use our *Base* model across several simulated scenarios. The ego is shown as a red rectangle, the replayed expert as the blue rectangle, and other (replayed) agents are yellow. The planned trajectory of the ego for the current time step is shown

Table 10: $\ell_2$ (with yaw) metric for each scenario tag.

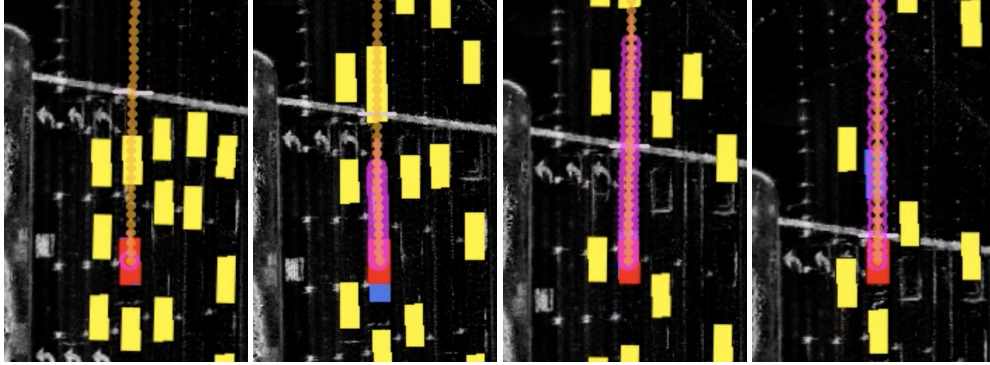| Tags | All | Straight | Right | Left | Stopped | Slow | Intersec. | Close | ASV |
|---|---|---|---|---|---|---|---|---|---|
| Base (ours) | 2.29 | 2.15 | **4.50** | **4.24** | 0.37 | 2.22 | **3.54** | **3.33** | 1.73 |
| Without focal | 2.39 | 2.24 | 5.76 | 5.09 | 0.50 | **2.12** | 3.69 | 3.52 | 1.74 |
| Balanced set + without focal | **2.28** | **2.12** | 5.16 | 5.07 | **0.34** | 2.12 | 3.67 | 3.54 | **1.57** |

Figure 8: The ego vehicle starts to move along side the vehicle ahead. Video clip: sim-start-from-stop1.mp4

in pink. In particular, Fig. 8 and Fig. 9 show maneuvers where the ego vehicle is standing still and starts to move along side the vehicle ahead. Fig. 10 and Fig. 11 show the ego vehicle slowing down and stopping behind a vehicle. Adaptive Cruise Control (ACC) scenarios where the ego adjusts its speed to follow the lead vehicle are shown in Fig. 12 and Fig. 13. Fig 14 and Fig. 15 depict cut-in scenarios where the ego vehicle reacts and plan a shorter trajectory to avoid collisions, moreover, the ego adjusts its speeds or comes to a complete stop depending on the behavior of the new vehicle ahead. Finally, Fig. 16 and Fig. 17 illustrate scenarios where the ego vehicle slows down to take a turn while making progress towards the goal.

Table 11: Simulation video clips

| Video clip | Snapshots | Duration [s] |
|---|---|---|
| **Starting from a stop** | | |
| sim-start-from-stop1.mp4 | Figure 8 | 10 |
| sim-start-from-stop2.mp4 | Figure 9 | 10 |
| **Stopping behind a vehicle** | | |
| sim-stop-for-vehicle1.mp4 | Figure 10 | 10 |
| sim-stop-for-vehicle2.mp4 | Figure 11 | 10 |
| **Adaptive cruise control** | | |
| sim-acc1.mp4 | Figure 12 | 10 |
| sim-acc2.mp4 | Figure 13 | 10 |
| **Moderate cut-ins** | | |
| sim-cut-ins1.mp4 | Figure 14 | 10 |
| sim-cut-ins2.mp4 | Figure 15 | 10 |
| **Turning** | | |
| sim-turn1.mp4 | Figure 16 | 10 |
| sim-turn2.mp4 | Figure 17 | 10 |

## A.9 Real-world driving video clips

Tab. 12 provides a curated list of real-world driving video clips collected during our testing in Las Vegas Strip and nearby areas. The description for each video clip is included in the corresponding snapshots figure caption.

Figure 9: Another case where the ego vehicle starts to move along side the vehicle ahead. Video clip: sim-start-from-stop2.mp4



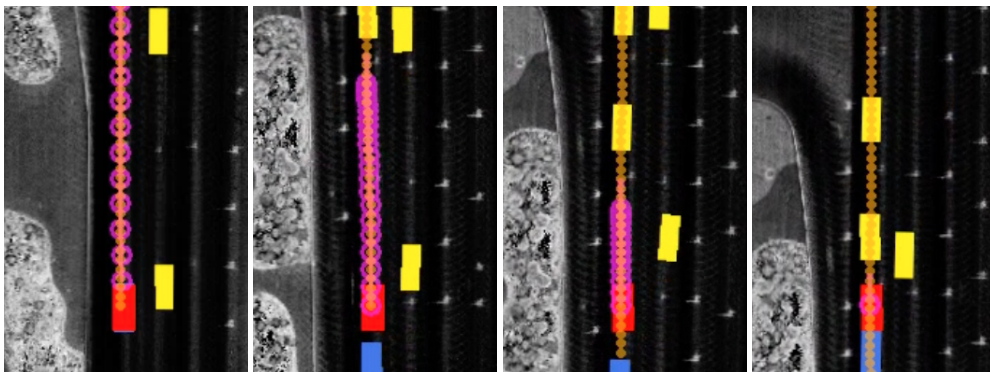Figure 10: The ego vehicle smoothly stops behind a vehicle. Video clip: sim-stop-for-vehicle1.mp4



Figure 11: The ego vehicle smoothly stops, from a higher velocity, behind a vehicle. Video clip: sim-stop-for-vehicle2.mp4
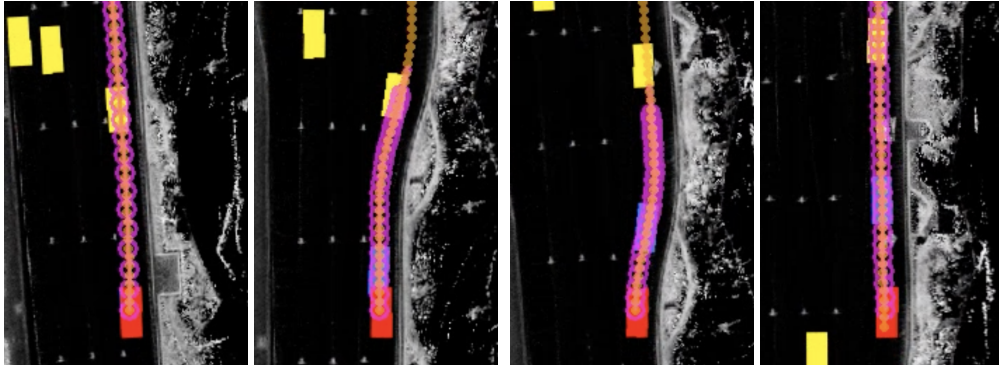
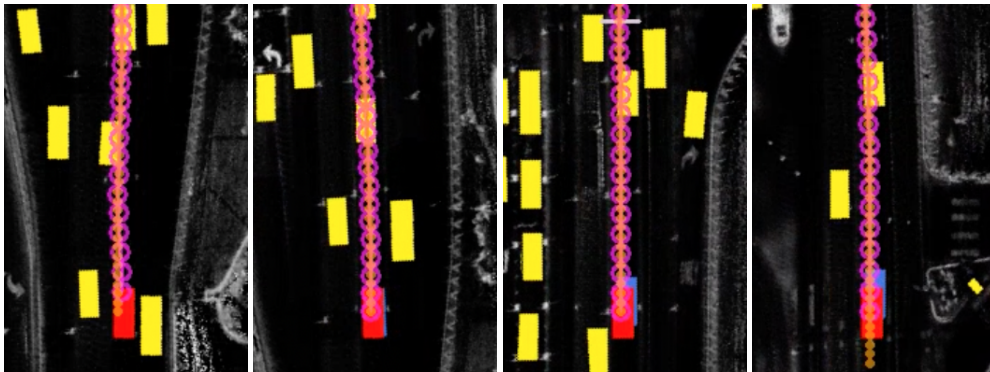Figure 12: The ego vehicle performs adaptive cruise control and slows down due to the behavior of the lead vehicle. Video clip: sim-acc1.mp4



Figure 13: The ego vehicle performs adaptive cruise control to maintain a safe distance from the lead vehicle. Video clip: sim-acc2.mp4



Figure 14: The ego vehicle handles a moderate cut-in maneuver. Video clip: sim-cut-ins1.mp4

Figure 15: The ego vehicle handles a moderate cut-in maneuver and stop behind the new lead vehicle. Video clip: sim-cut-ins2.mp4



Figure 16: The ego slows down and makes a right turn. Video clip: sim-turn1.mp4



Figure 17: The ego vehicle makes a 90° right turn while maintaining a comfortable speed. Video clip: sim-turn2.mp4

Table 12: Real-world driving video clips

| Video clip | Snapshots | Duration [s] |
|---|---|---|
| **Driving with a vehicle far ahead** | | |
| far-ahead-01.mp4 | Figure 18 | 34 |
| far-ahead-02.mp4 | Figure 19 | 25 |
| far-ahead-03.mp4 | Figure 20 | 20 |
| far-ahead-04.mp4 | Figure 21 | 30 |
| **Driving with a vehicle close ahead** | | |
| close-ahead-01.mp4 | Figure 22 | 40 |
| close-ahead-02.mp4 | Figure 23 | 30 |
| close-ahead-03.mp4 | Figure 24 | 30 |
| **Passenger pickup/dropoff zones** | | |
| pudo-01.mp4 | Figure 25 | 30 |
| pudo-02.mp4 | Figure 26 | 26 |
| **Stoping behind a vehicle** | | |
| stopping-01.mp4 | Figure 27 | 15 |
| stopping-02.mp4 | Figure 6 | 15 |
| stopping-03.mp4 | Figure 28 | 15 |
| **Cut-ins** | | |
| cut-in-01.mp4 | Figure 29 | 15 |
| cut-in-02.mp4 | Figure 30 | 20 |
| cut-in-03.mp4 | Figure 31 | 32 |



Figure 18: The ego vehicle follows dense traffic as it starts to flow after a transition from red to green traffic light. Video clip: far-ahead-01.mp4



Figure 19: The ego vehicle follows traffic as it starts to flow after a transition from red to green traffic light. Video clip: far-ahead-02.mp4



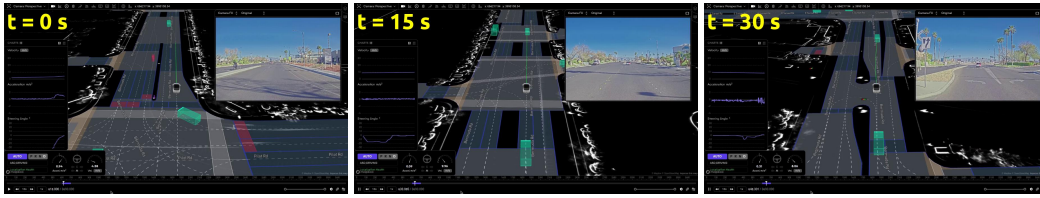Figure 20: The ego vehicle drives surrounded by few vehicles. Video clip: far-ahead-03.mp4

Figure 21: The ego vehicle performs a smooth lane-change. Video clip: far-ahead-04.mp4



Figure 22: The ego vehicle drives with a vehicle close ahead, performs a smooth lane-change and merges into a perpendicular road. Video clip: close-ahead-01.mp4



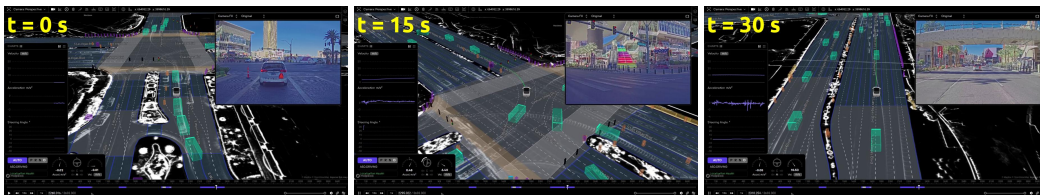Figure 23: The ego vehicle slows down and speeds up to match the lead vehicle speed. Video clip: close-ahead-02.mp4



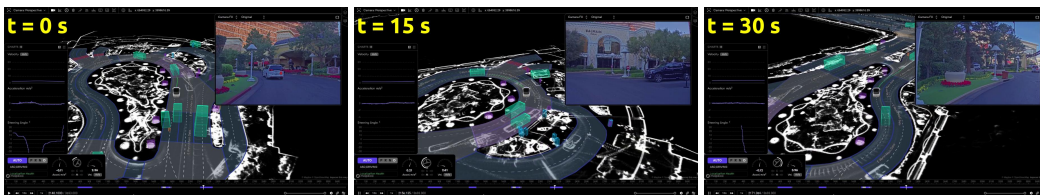Figure 24: The ego vehicle makes a protected left turn following the green traffic light. Video clip: close-ahead-03.mp4



Figure 25: The ego vehicle drives through a passenger pickup/dropoff zone. Video clip: pudo-01.mp4
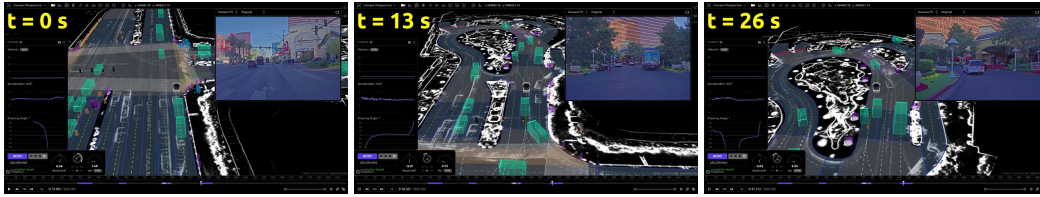
Figure 26: The ego vehicle makes an unprotected right turn to reach a passenger pickup/dropoff zone. Towards the end of the maneuver, the safety driver takes over to adjust the ego vehicle route. Video clip: pudo-02.mp4



Figure 27: The ego vehicle makes a hard stop from 10.5 m/s behind a vehicle. Towards the end of the maneuver, the safety driver takes over a comfortable stopping distance. Video clip: stopping-01.mp4
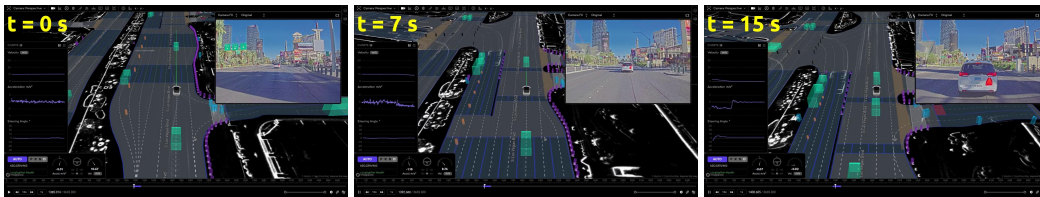


Figure 28: The ego vehicle smoothly stops from 10.5 m/s behind a vehicle. Video clip: stopping-03.mp4
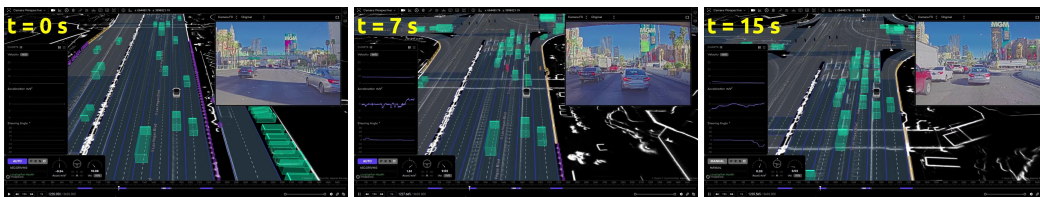


Figure 29: Edge case where our planner does not handle properly a sharp cut-in and the safety driver has to take over. Video clip: cut-in-01.mp4



Figure 30: The ego vehicle reacts properly to a moderate cut-in manuever at 10.5 m/s. Video clip: cut-in-02.mp4

Figure 31: The ego vehicle reacts properly to a sharp cut-in manuever at 10.5 m/s. Video clip: cut-in-03.mp4