

Learning to Drive from a World Model

Mitchell Goff Greg Hogan George Hotz Armand du Parc Locmaria Kacper Raczy
Harald Schäfer Adeeb Shihadeh Weixing Zhang Yassine Yousfi
comma.ai
autonomy@comma.ai

Abstract

Most self-driving systems rely on hand-coded perception outputs and engineered driving rules. Learning directly from human driving data with an end-to-end method can allow for a training architecture that is simpler and scales well with compute and data.

In this work, we propose an end-to-end training architecture that uses real driving data to train a driving policy in an on-policy simulator. We show two different methods of simulation, one with reprojective simulation and one with a learned world model. We show that both methods can be used to train a policy that learns driving behavior without any hand-coded driving rules. We evaluate the performance of these policies in a closed-loop simulation and when deployed in a real-world advanced driver-assistance system.

1. Introduction

Autonomous driving has seen remarkable progress in recent years, with learning-based approaches replacing increasing portions of the system. However, despite these advancements, most self-driving products still rely on handcrafted rules on top of a layer of perception. These modular methods require significant engineering effort to generalize to real-world complexities and edge cases. Instead, End-to-End (E2E) learning offers a more scalable solution by training a driving policy to imitate human driving behavior from real data. An E2E policy can take in raw sensor inputs, such as images, and directly output a driving plan or control action, eliminating the need for manual rule design [4].

A key challenge in E2E learning is how to train a policy that can perform well under the non-i.i.d. assumption made by most supervised learning algorithms such as Behavior Cloning [1]. In the real world, the policy’s predictions influence its future observations. Small errors accumulate over time, leading to a compounding effect that drives the system into states it never encountered during training.

To overcome this, the driving policy needs to be trained on-policy, allowing it to learn from its own interactions with

the environment, and enabling it to recover from its own mistakes. Running on-policy learning in the real world is costly and impractical [13], making simulation-based training essential.

Traditional driving simulators are often handcrafted with explicit limited traffic behaviors and scenes. While useful for testing, these simulators fail to capture the full complexity and richness of real-world driving.

In this work, we explore how two data driven simulators can be used to train an E2E driving policy: a reprojective novel view synthesis simulator [27], and a learned World Model [10, 11, 25]. By using real-world data these simulators can capture the full diversity of real-world driving scenarios, and provide ground-truth for policy decisions by imitating the human driving decisions. We propose a method to distill human driving behaviors during on-policy training, by anchoring a supervising model to future states.

We discuss limitations of the reprojective simulator, and the potential of the World Model simulator to scale with data and compute to overcome these limitations. We show that policies trained in this way learn normal driving behaviors, such as staying in a lane and changing lanes, and that they can be used in real-world applications, such as an Advanced Driver Assistance System (ADAS).

To our knowledge, this is the first work to show how end-to-end training, without handcrafted features, can be used in a real-world ADAS. Additionally, we believe this is the first use of a world model simulator for on-policy training of a policy that is deployed in the real world.

2. Formulation

2.1. Driving Policy

Our goal is to learn an End-to-End (E2E) driving policy π that maps from a history of observations (o_1, o_2, \dots, o_T) to a distribution over next actions a_{T+1} . We consider a history h_T^π to be a sequence of observations and previous actions.

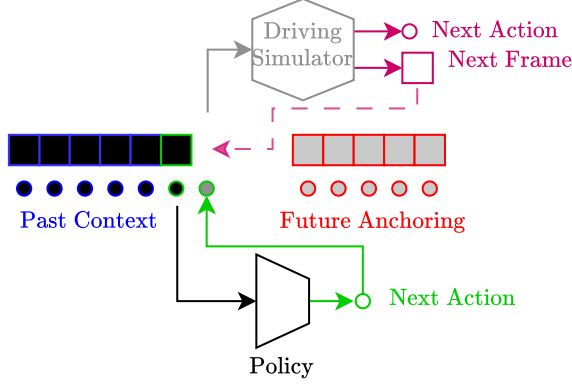


Figure 1. One step of the World Model Simulation rollout. Gray filled shapes are inputs to the World Model. Black filled shapes are inputs to both the Policy Model and the World Model (note that the Policy Model can be the World Model itself). Circles are actions (positions and orientations) and rectangles are observations (images).

$$\pi : h_T^\pi \mapsto p(a_{T+1} | h_T^\pi),$$

$$\text{with } h_T^\pi = \left((o_1, a_1), (o_2, a_2), \dots, (o_T, a_T) \right). \quad (1)$$

The action space \mathcal{A} is defined as a desired turning curvature and a desired longitudinal acceleration. The observation space \mathcal{O} is defined as a set of camera images only (Vision Only Policy). For simplicity, we will only show images coming from a single camera (narrow field of view). In practice, we use images from two cameras: wide and narrow field of view in order to capture a larger portion of the scene. The formulation can be extended to include more cameras without loss of generality.

We are given a dataset of expert demonstrations $\mathcal{D} = \left\{ \left((s_1, a_1), \dots, (s_T, a_T) \right) \right\}_{i=1}^n$. We aim to learn a driving policy π given the expert demonstrations in \mathcal{D} .

The state space \mathcal{S} is defined as the set of camera images, but can also include other sensor data such as GPS, and IMU data. In this work, we restrict the state space to the set of camera images \mathcal{O} and a global pose estimate (position and orientation) of the vehicle $p_t = (x, y, z, \phi, \theta, \psi) \in \mathcal{P} \subset \mathbb{R}^6$.

The global pose is obtained using a tightly coupled GPS/Vision Multi-State Constraint Kalman Filter (MSCKF) [14, 17, 26].

2.2. Driving Simulator

We define a Driving Simulator as composed of (i) a Driving State Generator and (ii) an Action Ground Truth Source.

The Driving State Generation can be based on traditional driving simulators such as CARLA [6], MetaDrive [15], etc.

or based on a so-called World Model that is learned from data (Section 4), or based on reprojective novel view synthesis techniques (Section 3).

The Action Ground Truth Source provides the supervision signal for training the driving policy. We describe a data-driven Action Ground Truth Source in Section 2.6.

2.3. Vehicle Model

A vehicle’s motion is described as a sequence of poses (p_1, p_2, \dots, p_T) . To simulate the effects of actions taken in simulation, a function is needed that produces poses based on actions. We call this a Vehicle Model [12]. Our Vehicle Model is designed to model a variety of real-world effects including vehicle dynamics, delayed steering response, wind, and more. Simulating these effects is needed for transferring policies trained in simulation to the real world, and is often referred to as Domain Randomization (sim2real) [20, 24, 28]. By inverting the Vehicle Model we can also estimate actions needed to achieve a trajectory of poses.

Figure 1 illustrates the different building blocks involved in one step of a driving simulator.

2.4. World Model

We define a World Model w as a stochastic model that predicts a future state given a history of states and actions. In order to make the World Model independent of the Vehicle Model described in 2.3, we consider the actions (desired curvature and acceleration) to be implicitly deducible from the vehicle’s poses (p_1, p_2, \dots, p_T) and the Vehicle Model. In other words, w maps a history of images and poses and next pose to a distribution of the next image.

$$w : h_T^w \mapsto p(o_T | h_T^w),$$

$$\text{with } h_T^w = \left((p_1, o_1), (p_2, o_2) \dots, (p_T,) \right). \quad (2)$$

Note that using the pose as the transition signal for the World Model enables augmenting the Vehicle Model’s parameters without needing to retrain the World Model.

2.5. Future Anchored World Model

We can train non-causal World Models similar to [2] conditioned on future observations and actions parametrized by $F = (f_s, f_e)$, where f_s is the start of the future horizon and f_e is the end of the future horizon. With $f_s > T$. This model can only be used offline, but has the advantage of predicting human-like driving video sequences and trajectories that converge to a goal state at F . We refer to this as recovery pressure.

$$\begin{aligned}
w : h_{T,F}^w &\mapsto p(o_T \mid h_{T,F}^w), \\
\text{with } h_{T,F}^w &= \left((p_{f_s}, o_{f_s}), \dots, (p_{f_e}, o_{f_e}), \right. \\
&\quad \left. (p_1, o_1), (p_2, o_2), \dots, (p_T,) \right). \tag{3}
\end{aligned}$$

2.6. Driving Ground Truth

The Action Ground Truth at time T refers to the actions $a_T \mid h_T$ the policy should take, given the past observations and the past actions it has taken, resulting in good driving behavior.

To generate this ground truth, we train the Future Anchored World Model to also predict the next pose or trajectory of poses \mathcal{T} given a history of observations and poses and a Future Anchoring, $\mathcal{T} \mid h_{T,F}$. When running in this mode, the World Model is referred to as a Plan Model, and the predicted trajectory can be mapped to ground-truth actions using the Vehicle Model.

Future Anchoring is essential for enabling the Plan Model to produce a trajectory that converges to a desirable goal state, F , regardless of the current state of the simulation, without it, the Plan Model does not exhibit recovery pressure when in a bad state.

Note that the Plan Model can be a separate model, but it is often trained jointly with the World Model to leverage shared representations. The Plan Model can also be used independently of the state generation method used in the simulator, i.e. it can be used with a reprojective simulation or a learned World Model.

3. Reprojective Simulation

Given a dense depth map d_T , a pose p_T , and an image o_T , we can render a new image o'_T by reprojecting the 3D points in the depth map to the new pose p'_T . This process is called Reprojective Simulation. In practice, we can use a history of images and depth maps to reproject the image and inpaint the missing regions. An example is shown in Figure 2.

3.1. Limitations of Reprojective Simulation

We list some of the limitations of Reprojective Simulation:

Assumption of a static scene: This formulation assumes that the scene is static, and does not depend on p_T , which is usually not the case. For example, swerving towards a neighboring car might cause the driver of the neighboring car to react. We refer to this issue as the counterfactual problem.

Depth estimation inaccuracies: The depth map d_T is usually noisy and inaccurate, which leads to artifacts in the reprojected image o'_T .

Occlusions: Regions that are occluded in o_T ought to be inpainted in o'_T , which is a challenging task, and also leads to artifacts in the reprojected image.

Reflections and lighting: By definition, Reprojective Simulation ignores the physics of light transport. Without ray tracing or an equivalent lighting model, reprojection can only tell where surfaces are, but not how light interacts with those surfaces from a new angle. This is a major limitation for night driving scenes, leading to noticeable lighting artifacts in the reprojected image, as shown in Figure 3.

Limited Range: The more p'_T differs from p_T , the more pronounced the artifacts become. In order to limit the artifacts, we need to limit the range of simulation to small values (typically less than 4m in translation), which is especially limiting for longitudinal motion.

Artifacts are correlated with $p'_T - p_T$: The artifacts in the reprojected image are correlated with the difference between the poses p'_T and p_T . This correlation is exploited by the policy to predict the future action, which is not desirable. We refer to this as cheating or shortcut learning [9].

4. World Models Simulation

The World Model simulator is a Future Anchored World Model and Plan Model.

4.1. Video Encoder

We use the pretrained Stable Diffusion image VAE [23]. More specifically, `vae-ft-mse-840000-ema-pruned` which has a compression factor of 8×8 and 4 latent channels per image. For simplicity, we exchangeably use o for the latent representation of the camera image from the VAE tokenizer in this section.

4.2. Diffusion Transformer

4.2.1. Architecture

We use the Diffusion Transformer (DiT) architecture [19], adapted to 3 dimensional inputs by extending the input/output patching table to a 3 dimensional table, then flattening all 3 dimensions before the Transformer blocks.

Similar to [3], the vehicle poses, world timesteps, and diffusion noise timesteps are used as conditioning signals for the Diffusion Transformer. The conditioning signal embeddings are summed and passed to the Adaptive Layer Norm layer (AdaLN) [32]. The AdaLN layer is modified to support different conditioning vectors along the time dimension.

Additionally, the attention layers use a block-wise (frame-wise) triangular causal mask, i.e. query tokens can only attend to (key, value) tokens within the same frame or the previous frames in the sequence. Note that this does not make the model physically causal, as future observations and future poses are prepended to the input sequence. However, this masking is required to use key-value caching (kv-caching) during inference, which is essential for efficient sampling.

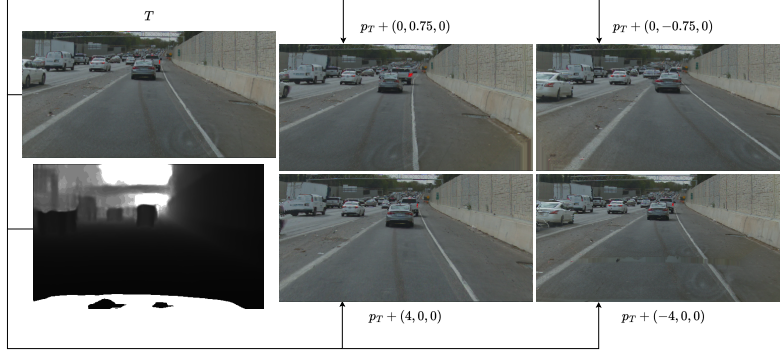


Figure 2. Left: Top: Image at T . Bottom: Depth map at T . Right: Reprojected images at T using 4 different translation vectors.

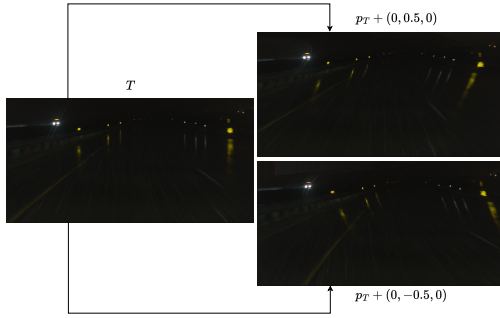


Figure 3. Left: Image at T , Right: Reprojected Images at T . Notice the lighting artifacts in the reprojected images.

To make it a Plan Model, the Transformer is equipped with a Plan Head, which is a stack of residual Feed Forward blocks. The Plan Head predicts the trajectory \mathcal{T} .

4.2.2. Training objective

We adopt the Rectified Flow (RF) objective [16] for training the Conditional Diffusion Transformer. For simplicity, we omit the subscripts T and F for the world timestep in the following equations. We sample the noise timestep $\tau \sim \text{Logit-Normal}(0.0, 1.0)$ [8] and noise the observations o using Equation 4.

$$o_\tau = \tau \epsilon + (1 - \tau) o \quad (4)$$

The Plan Head output \mathcal{T} uses a Multi-hypothesis Planning loss (MHP) [5] with 5 hypotheses. Each hypothesis is trained using a heteroscedastic Negative Log Likelihood (NLL) loss with a Laplace prior [18].

The total loss \mathcal{L} is a weighted sum of the Rectified Flow loss \mathcal{L}_{RF} and the MHP loss $\mathcal{L}_{\mathcal{T}}$ with a hyperparameter α as described in Equation 5.

$$\mathcal{L} = \mathcal{L}_{\text{RF}} + \alpha \mathcal{L}_{\mathcal{T}}$$

$$\text{where } \mathcal{L}_{\text{RF}}(o, p, \epsilon, \tau) = \|w(o_\tau, p, \tau) - (o - \epsilon)\|^2 \quad (5)$$

$$\mathcal{L}_{\mathcal{T}}(o, p, \epsilon, \tau) = \text{MHP}(w(o_\tau, p, \tau), \mathcal{T})$$

4.3. Sequential Sampling

At every world timestep T we use a simple Euler discretization with 15 steps $\Delta\tau = 1/15$ to sample the next latents \tilde{o}_T , the sampling process follows equation 6. None of the context latents ($o_{f_s}, \dots, o_{f_e}, o_1, \dots, o_{T-1}$) are noised, and we use $\tau = 0$ as input to the model for those timesteps. The vehicle position p_T can be sampled from the World Model’s own Plan Head, from a Policy, from the ground truth trajectory, or artificially crafted, as described in Section 4.6.

$$\tilde{o}_{\tau+\Delta\tau} = \Delta\tau w(\tilde{o}_\tau, p, \tau + \Delta\tau) + \tilde{o}_\tau \quad (6)$$

For the next timestep $T + 1$ we shift the context latents by one timestep, and append the sampled latent \tilde{o}_T to the context latents ($o_{f_s}, \dots, o_{f_e}, o_2, \dots, o_{T-1}, \tilde{o}_T$). We repeat this process until we reach the future horizon f_e .

4.4. Noise Level Augmentation

In order to make the model robust to the so-called ”auto-regressive drift” [29] i.e. errors in the Sequential Sampling process compounding frame by frame, we use a noise level augmentation technique.

For some training samples (with probability $p = 0.3$), we do the following:

- Sample different noise levels $\tau \sim \text{Logit-Normal}(0.0, 0.25)$ at world timesteps $(1, \dots, T - 1)$,
- Don’t noisify the future anchoring latents, i.e. $\tau = 0$ at world timesteps (f_s, \dots, f_e) ,
- Input $\tau = 0$ to the model at world timesteps $f_s, \dots, f_e, 1, \dots, T - 1$,
- Only compute the diffusion loss on the latents at T .

This diffusion noise level augmentation was essential to making the model robust to accumulated errors in the Sequential Sampling process. A similar technique was proposed in [29], although we didn’t find it necessary to discretize the noise levels.



Figure 4. Five examples of World Model simulation. Blue bordered frames are the last frames of the past context, red bordered frames are the first frames of the future anchoring, and green bordered frames are simulated frames. Notice how the simulated frames comply with the future anchoring by executing lanes changes, or turning the traffic light to green.

4.5. Implementation Details and Data

We train three sizes of DiTs based on configurations from the GPT-2 models [21] `gpt` (250M parameters), `gpt-medium` (500M parameters), and `gpt-large` (1B parameters). We use three dataset sizes: 100k, 200k, and 400k segments, each segment is 1 minute long of driving video and vehicle poses. The videos are downsampled to 128×256 pixels before being fed to the VAE. We down-sample all the data to 5 Hz.

Data and Model size scaling results are shown in Figure 5. Unless otherwise stated, we use the 500M DiT model trained on 400k segments for the rest of the experiments.

Every training sample is constructed as follows: we sample a context of $T = 2s$ from the dataset, then we sample a world timestep $T < f_s < 9s$: the start of the future horizon, $f_e - f_s$ is kept constant to 1s. The image VAE features and the vehicle poses are then concatenated as described in equation 3.

At every timestep, the trajectory \mathcal{T} is constructed as a sequence of positions, speeds, accelerations, orientations, and orientation rates up to a future horizon of 10 seconds.

See Figure 4 for examples of World Model rollouts.

4.6. World Model Evaluation

We use LPIPS [33] similarity of the generated images to the ground truth images as a measure of image/video quality.

Note that the baseline LPIPS score is $\text{LPIPS} = 0.148$ due to the VAE compression as measured on the test set, and the lower the LPIPS score the better the quality of the generated images.

In order to evaluate how accurately the World Model respects the vehicle positions p_T inputs, we use a Pose Net to

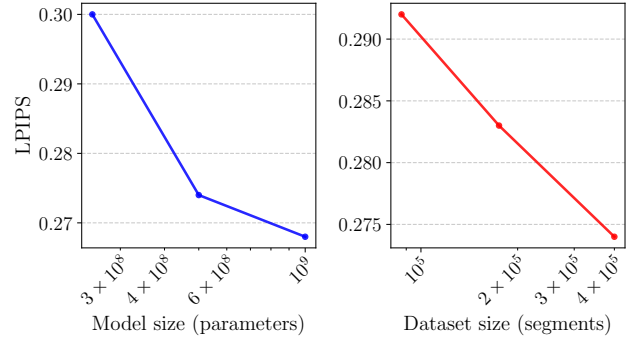


Figure 5. Left: LPIPS for different DiT model sizes, trained on 400k segments. Right: LPIPS for different dataset sizes, for a DiT of 500M parameters. Both are from the action teacher-forced sequential rollout setting.

measure the error between the commanded vehicle motion and the one generated by the World Model.

The Pose Net is a supervised model trained to predict a variety of outputs, such as pose, lane lines, road edges, lead car position, etc.

We evaluate the World Model over 1,500 rollouts from different segments of the test set. We can run the World Model in multiple modes.

4.6.1. Image Quality

Observation and action teacher-forced next frame prediction:

At each world timestep T , we sample a frame \tilde{o}_T using actions from the ground truth trajectory (rather than those from a policy). Then we replace it with the ground truth image o_T before placing it in the context latents for the next timestep $T + 1$. Results of this test are shown in left Figure 6.

MAE	non VAE Compressed	VAE Compressed
x speed	0.46366 m/s	0.59390 m/s
y speed	0.04216 m/s	0.04393 m/s
z speed	0.04424 m/s	0.04548 m/s
roll rate	0.00468 rad/s	0.00524 rad/s
pitch rate	0.00433 rad/s	0.00453 rad/s
yaw rate	0.00211 rad/s	0.00254 rad/s
y lane lines	0.15852 m	0.15995 m

Table 1. Pose Net MAE on non-VAE compressed and VAE compressed segments.

4.6.2. Video Quality

Action teacher-forced sequential rollout: Here we sample a frame \tilde{o}_T using actions from the ground truth trajectory, and place it in the context latents for the next timestep $T+1$. This mode is equivalent to using the World Model normally, but with ground truth actions rather than those from a policy. Results of this test are shown in right Figure 6.

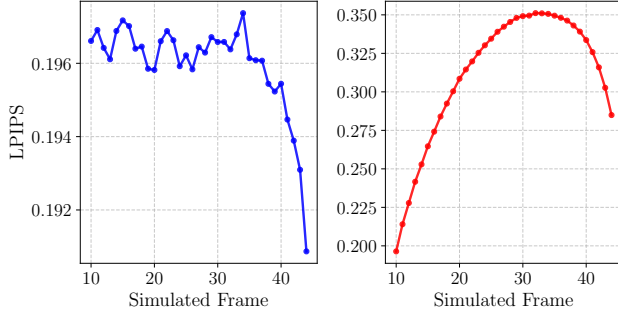


Figure 6. LPIPS for right: observation and action teacher-forced next frame prediction (image quality evaluation), left: action teacher-forced sequential rollout (video quality evaluation).

4.6.3. Pose accuracy

Action policy-forced sequential rollout: Here we sample the next frame \tilde{o}_T using actions from a policy π . This is the general mode of operation for the World Model. The policy can also be a so-called "noise model". One example of a noise model is a policy that deviates from the starting position by a fixed distance as demonstrated in Figure 8. Note that when the policy is the ground truth trajectory, this is equivalent to the Action teacher-forced sequential rollout.

The pose errors of the World Model measured by the Pose Net are shown in Figure 7.

We also run the World Model with the following noise model, we force a smooth lateral deviation of $\pm 0.5\text{m}$ over the first 25 steps of the rollout, then let the World Model recover. We show an example of this noise in Figure 8. Figure 9 shows the commanded lateral deviation, and the

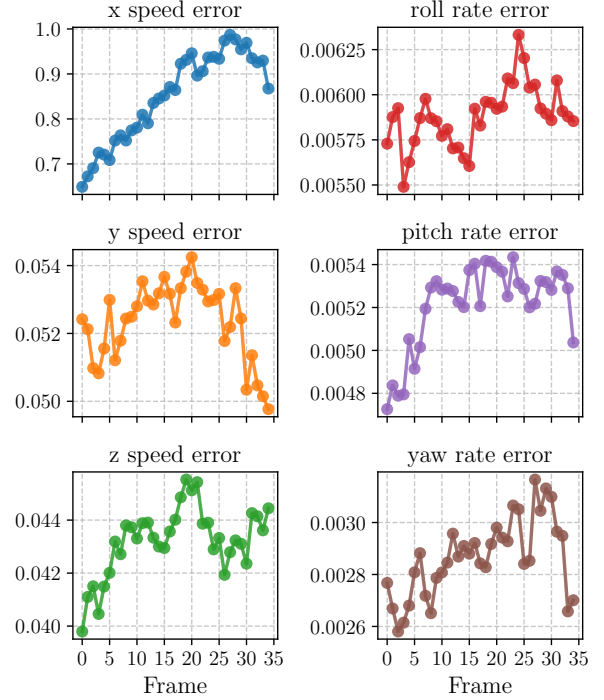


Figure 7. Action teacher-forced sequential rollout pose errors.

actual deviation (measured by the Pose Net) simulated by the World Model, averaged over 1,500 rollouts. The World Model simulates the commanded deviation, but not to its full extent.

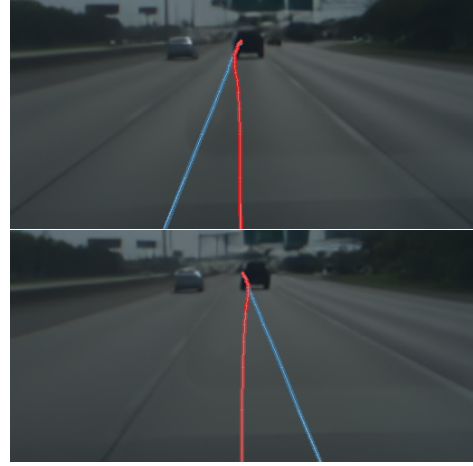


Figure 8. Action noise-forced sequential rollout. We force a smooth lateral deviation of $\pm 0.5\text{m}$ over the 25 steps of simulation. Pictured above is the deviation to the right (top) and to the left (bottom) at step 25.

5. Driving Policy Training

The driving policy π is a stack of two Neural Networks.

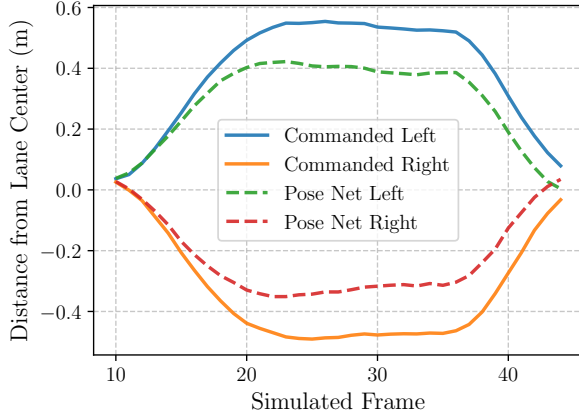


Figure 9. Deviation from the lane center in noise forced simulation. Dashed lines indicate deviation measured by the Pose Net, and solid lines indicate the commanded deviation.

The first is a supervised feature extractor based on the FastViT architecture [30], which is trained to predict a variety of outputs including lane lines, road edges, lead car information, and ego car future trajectory. Note that lane lines and road edges outputs are used for visualization, and never used as part of a steering policy.

The second is a small Transformer [31] based temporal model predicting the same outputs as the feature extractor, in addition to the next action (desired curvature and acceleration). The temporal model’s inputs are the features from the (frozen) FastViT extractor over the last 2 seconds.

Similar to the plan head of the World Model, the trajectory output of the driving policy is trained using MHP loss with 5 hypotheses and a Laplace prior. The other outputs are trained using NLL loss with a Laplace prior.

We distinguish two different approaches to training the temporal model part of the policy. Off-Policy Learning and On-Policy Learning. Off-policy learning refers to using supervised learning on the collected dataset of expert demonstrations. This is also known as imitation learning. We describe On-Policy Learning in Section 5.1.

5.1. On-Policy Learning

The temporal model is trained using a driving simulator, the policy’s training samples come from its own interaction with the environment. We adopt a similar architecture to IMPALA [7] where a set of actors running in parallel generate experiences that are sent to a central learner. The learner updates the policy then sends a new version to the actors using a parameter server.

These experiences (also referred to as rollouts) are se-

quences of the form:

$$h^{\pi, wp} = \left((o_1, a_1, \hat{a}_1^{wp}), (o_2, a_2, \hat{a}_2^{wp}), \dots, (o_{f_s}, a_{f_s}, \hat{a}_{f_s}^{wp}) \right), \quad (7)$$

where \hat{a}^{wp} is the action derived from the (plan equipped) Future Anchored World Model wp during the rollout.

At every timestep T , the actor takes the action a_T sampled from the latest policy π available in the parameter server at the start of the rollout. The driving simulator generates the observations o_T given the state of the world and the commanded action. The simulator can be wp itself, a different World Model w , or any driving simulator.

The rollouts end when we reach the future horizon f_s , after which they are sent to the learner.

The learner optimizes the mapping $\pi : h_T^\pi \mapsto p(\hat{a}_T^w | h_T^\pi)$, i.e. the policy learns to predict the actions that the World Model would take given the history h_T^π .

To ensure the policy is robust to a variety of real-world effects, the Vehicle Model described in 2.3 must model a complex distribution of action responses during training.

5.2. Information Bottleneck

To prevent the policy from exploiting simulator-specific artifacts described in Section 3.1, we regularize the feature extractor by limiting the amount of information it can output to roughly 700bits. We impose this limit by adding white Gaussian noise during training, the bottleneck can be interpreted as a Gaussian communication channel with a per-sample information capacity: $\frac{1}{2} \log(1 + \text{SNR})$. This bottleneck is similar to Gaussian Dropout [22] which uses multiplicative noise instead of additive noise.

5.3. Policy Evaluation Suite

We focus on evaluating the policy’s lateral driving performance, i.e. its ability to accurately and smoothly steer to maintain human-like lane positioning, and successfully execute lane-changes under various conditions. Without loss of generality, longitudinal metrics and tests can be similarly integrated and are the subject of future work.

5.3.1. Simulated On-Policy Unit Tests

We use the MetaDrive Simulator [15] to evaluate the policy in closed loop. See Figure 10 for rendered examples. We define a set of unit tests that the policy should pass:

Convergence to lane-center test on straights and turns (24 scenarios): We test whether the policy converges to a good position in the lane regardless of small offsets in its starting conditions. Note that the position in the lane doesn’t need to be the centered, but rather a position that a human driver would converge to. See left Figure 11.

Lane change completion test (20 scenarios): We test whether the policy can complete a lane change maneuver

regardless of small offsets in its starting conditions in the lane. During training a conditioning impulse is added prior to lane changes, and we input the same impulse during inference to trigger a lane change. See right Figure 11.



Figure 10. Examples of the MetaDrive simulated unit test scenarios.

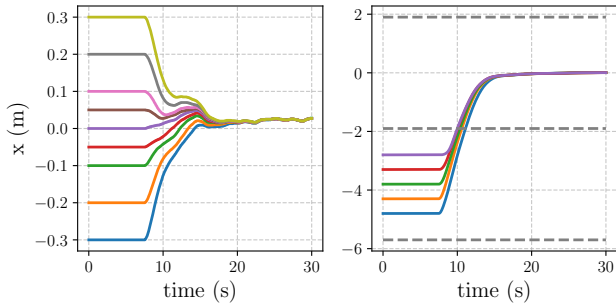


Figure 11. Left: Convergence to lane-center test results. Right: Lane change completion test results.

5.3.2. Off Policy Evaluation

We use a holdout set of 1,500 segments from the dataset. We evaluate the policy on these segments off-policy by running it at each timestep and computing a variety of metrics. For simplicity, we only report the trajectory overall MAE.

5.3.3. In the field Evaluation

We evaluate the policy in the field by deploying it to openpilot¹ and collecting data from users. openpilot is an open-source ADAS which supports a wide variety of production vehicles. The end-to-end policies described here directly control the steering actions of the vehicle to provide continuous auto-steering when the system is engaged. The longitudinal action is controlled by a classical ACC (Adaptive Cruise Control) policy, that uses lead detection and radar to slow down for other vehicles and maintains cruise speed.

5.4. Results

We evaluate three different policies: one trained off-policy, one trained on-policy using a reprojective simulator, and one trained on-policy using the World Model simulator. The results are shown in Table 2. We clearly see that the policy trained off-policy fails in the on-policy tests, despite performing better in the off-policy accuracy evaluation.

¹<https://github.com/commaai/openpilot>

	Off-policy	Reprojective	World Model
MetaDrive lane center	5/24	24/24	24/24
MetaDrive lane change	8/20	20/20	19/20
Off-policy trajectory MAE	0.361	0.369	0.394

Table 2. Performance of driving policies trained in different conditions on the MetaDrive simulator and off-policy accuracy evaluation.

	Reprojective	World Model
Number of trips	47,047	40,026
Engaged % (time)	27.63%	29.92%
Engaged % (distance)	48.10%	52.49%

Table 3. Performance in the field of the trained driving policies.

Both on-policy learning methods have been successfully deployed in the real world in openpilot. Table 3 presents usage metrics collected over approximately two months of driving from a cohort of 500 users. Since openpilot is a level 2 system where disengagements are expected, we use the percentage of time and distance during which the system was engaged as the primary metric. These results demonstrate that both policies are capable of delivering meaningful driver assistance in real-world conditions.

6. Conclusion and Future Work

In this work, we propose an architecture for training a driving policy on-policy in a data-driven simulator based on real human driving data. This simulator produces both input images and action ground-truth based on real human driving data. We propose two different simulation strategies, one using a more traditional reprojective simulator, and another using a world model.

We show that the driving policies produced by these training strategies learn basic driving behaviors in our test suite, without any engineered behaviors during training. We also show that these policies can be used to produce useful ADAS products that assist driving in the real world.

We discuss the limitations of the reprojective simulation and how we expect the world models strategy will continue to scale with data and compute to produce better driving policies. While this work focuses on lateral driving policy, all methods generalize to longitudinal policy, in future work we expect to demonstrate useful ADAS products using end-to-end longitudinal policy trained with these methods.

References

- [1] Michael Bain and Claude Sammut. A framework for behavioural cloning. In *Machine intelligence 15*, pages 103–129, 1995. 1
- [2] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampe-dro, and Jeff Clune. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems*, 35:24639–24654, 2022. 2
- [3] Amir Bar, Gaoyue Zhou, Danny Tran, Trevor Darrell, and Yann LeCun. Navigation world models, 2024. 3
- [4] 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. 1
- [5] Henggang Cui, Vladan Radosavljevic, Fang-Chieh Chou, Tsung-Han Lin, Thi Nguyen, Tzu-Kuo Huang, Jeff Schneider, and Nemanja Djuric. Multimodal trajectory predictions for autonomous driving using deep convolutional networks. In *2019 international conference on robotics and automation (icra)*, pages 2090–2096. IEEE, 2019. 4
- [6] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. CARLA: An open urban driving simulator. In *Proceedings of the 1st Annual Conference on Robot Learning*, pages 1–16, 2017. 2
- [7] Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *International conference on machine learning*, pages 1407–1416. PMLR, 2018. 7
- [8] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. In *Proceedings of the 41st International Conference on Machine Learning*. JMLR.org, 2024. 4
- [9] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020. 3
- [10] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *Advances in Neural Information Processing Systems 31*, pages 2451–2463. Curran Associates, Inc., 2018. <https://worldmodels.github.io>. 1
- [11] Anthony Hu, Lloyd Russell, Hudson Yeo, Zak Murez, George Fedoseev, Alex Kendall, Jamie Shotton, and Gianluca Corrado. Gaia-1: A generative world model for autonomous driving, 2023. 1
- [12] Reza N Jazar. *Vehicle dynamics*. Springer, 2008. 2
- [13] 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. 1
- [14] Mingyang Li and Anastasios I. Mourikis. High-precision, consistent ekf-based visual-inertial odometry. *The International Journal of Robotics Research*, 32(6):690–711, 2013. 2
- [15] Quanyi Li, Zhenghao Peng, Lan Feng, Qihang Zhang, Zhenghai Xue, and Bolei Zhou. Metadrive: Composing diverse driving scenarios for generalizable reinforcement learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 2, 7
- [16] Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022. 4
- [17] Anastasios I. Mourikis and Stergios I. Roumeliotis. A multi-state constraint kalman filter for vision-aided inertial navigation. In *Proceedings 2007 IEEE International Conference on Robotics and Automation*, pages 3565–3572, 2007. 2
- [18] David A Nix and Andreas S Weigend. Estimating the mean and variance of the target probability distribution. In *Proceedings of 1994 IEEE international conference on neural networks (ICNN'94)*, pages 55–60. IEEE, 1994. 4
- [19] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 3
- [20] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3803–3810, 2018. 2
- [21] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 5
- [22] Mélanie Rey and Andriy Mnih. Gaussian dropout as an information bottleneck layer. In *NeurIPS Workshop on Bayesian Deep Learning*, 2021. 7
- [23] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 3
- [24] Fereshteh Sadeghi and Sergey Levine. CAD2RL: real single-image flight without a single real image. In *Robotics: Science and Systems XIII, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA, July 12-16, 2017*, 2017. 2
- [25] Eder Santana and George Hotz. Learning a driving simulator, 2016. 1
- [26] Harald Schafer, Eder Santana, Andrew Haden, and Riccardo Biasini. A commute in data: The comma2k19 dataset, 2018. 2
- [27] Steven M Seitz and Charles R Dyer. View morphing. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 21–30, 1996. 1
- [28] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization

for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 23–30, 2017. [2](#)

- [29] Dani Valevski, Yaniv Leviathan, Moab Arar, and Shlomi Fruchter. Diffusion models are real-time game engines. In *The Thirteenth International Conference on Learning Representations*, 2025. [4](#)
- [30] Pavan Kumar Anasosalu Vasu, James Gabriel, Jeff Zhu, Oncel Tuzel, and Anurag Ranjan. Fastvit: A fast hybrid vision transformer using structural reparameterization. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5785–5795, 2023. [7](#)
- [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. [7](#)
- [32] Jingjing Xu, Xu Sun, Zhiyuan Zhang, Guangxiang Zhao, and Junyang Lin. Understanding and improving layer normalization. *Advances in neural information processing systems*, 32, 2019. [3](#)
- [33] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. [5](#)