# Learning to Drive by Watching YouTube Videos: Action-Conditioned Contrastive Policy Pretraining Supplementary Material

Qihang Zhang<sup>1</sup>, Zhenghao Peng<sup>1</sup>, and Bolei Zhou<sup>2</sup>

 $^1\,$  The Chinese University of Hong Kong, Hong Kong SAR, China $^2\,$  University of California, Los Angeles, USA

# A YouTube Driving Dataset

Our YouTube Driving Dataset contains a massive amount of real-world driving frames with various conditions, from different weather, different regions, to diverse scene types. We collect 134 videos with a total length of over 120 hours, covering up to 68 cities. Figure 7 illustrates the videos' lengths distribution. For pretraining, we sample one frame every one second, resulting to a dataset of 1.30 million frames in total.



Fig. 7. Video length distribution

#### 2 Q. Zhang, Z. Peng, B. Zhou

Methods	IL Demonstration Size $(\times 40K)$			
niotiiotab	10%	20%	40%	100%
Random-f	$0.0{\pm}0.0$	$0.0{\pm}0.0$	$2.0{\pm}0.0$	$4.0{\pm}0.0$
ImageNet-f	$4.0 {\pm} 0.0$	$4.0{\pm}2.8$	$5.3 {\pm} 0.9$	$6.7{\pm}0.9$
MoCo-f	$10.7{\pm}2.5$	$30.7 {\pm} 4.1$	$38.7 {\pm} 1.9$	$54.7{\pm}6.6$
ACO-f	<b>24.0</b> ±1.6	<b>40.7</b> ±3.4	$\textbf{47.3}{\pm}0.9$	<b>61.3</b> ±3.8

Table 5. Success Rate of Imitation Learning with frozen backbone

 Table 6. Effects of Batch Normalization in Imitation Learning. "bn" stands

 for Batch Normalization

Methods	IL Demonstration Size $(\times 40K)$			
mothods	10%	20%	40%	100%
MoCo w/o bn	17.3±1.9	$55.3 \pm 16.8$	$59.3 \pm 5.7$	80.0±1.6
MoCo w/ bn	<b>19.3</b> ±3.4	<b>60.3</b> $\pm 8.2$	<b>76.0</b> $\pm 9.1$	80.7±2.5
ACO w/o bn	27.3±3.8	$34.0 \pm 4.3$	72.0±4.3	78.0±4.3
ACO w/ bn	<b>30.7</b> ±3.4	<b>66.0</b> $\pm 5.7$	<b>82.0</b> ±5.0	<b>96.0</b> ±3.3

# **B** Imitation Learning with frozen backbone

As shown in Table 5, we conduct imitation learning with frozen backbone. With backbone frozen, ACO demonstrates highest success rate. However, all methods suffer from severe performance drop compared to finetuned IL. For ImageNet and random pretrain, they drop by over 80% in full dataset size.

In contrast, RL methods exhibit less performance drop with frozen backbone. As shown in Fig. 5 in main paper, the biggest drop by frozen RL compared to finetuned RL is less than 50%. This phenomenon shows that freezing the backbone exacerbates the need for a good pretrain for IL. Policy trained with bad pretrain weight will give suboptimal actions, for example, driving the car to the road edge. IL, notoriously suffering from *distribution shift*, has no means to recover this error since the collected demonstrations do not contain any recovery action. Nevertheless, RL frequently visits dangerous states during exploration, and recovery strategies can hence be learned.

# C Effects of Batch Normalization on ACO

As suggested in [3], Batch Normalization could help alleviate the problem caused by different distributions of contrastive-learning-based parameters and supervised-learning-based parameters. In this section, we compare the imitation learning performance on ACO and MoCo with and without batch normalization in downstream IL's fine-tuning. As shown in Table 6, Batch Normalization is essential for both MoCo and ACO. Without adjusting parameters' distribution by batch normalization, these two contrastive learning methods show declining performance.



(B) Carla

Fig. 8. Reconstructed scenes by AutoEncoder in YouTube and Carla domain

# D Visualization of Auto-Encoder Baseline

To understand the poor performance of auto-encoder in all downstream tasks, we visualize the input and reconstructed image in both YouTube (training) and Carla (testing) domain in Figure 8.

It is obvious that much attention is paid to elements that are not relevant to driving decision-making in reconstructed images, such as buildings and the sky. On the contrary, the policy-relevant cues such as the lane lines and road edges are oversimplified. This means that useless information is learned in auto-encoder's feature and thus damages the performance of downstream tasks.

# **E** Implementation Details and Hyper Parameters

We list the hyper-parameters used in ACO pretraining in Table 7, inverse dynamics model training in Table 8, imitation learning in Table 9, and reinforcement learning in Table 10.

For imitation learning, following CILRS [1], the training target is to optimize L1loss between predicted value and ground truth on steering, throttle, brake, and velocity:

$$\mathcal{L}_{il} = \lambda_s \mathcal{L}_s + \lambda_t \mathcal{L}_t + \lambda_b \mathcal{L}_b + \lambda_v \mathcal{L}_v, \tag{1}$$

where  $\mathcal{L}_s$ ,  $\mathcal{L}_t$ ,  $\mathcal{L}_b$ , and  $\mathcal{L}_v$  stands for L1loss on steering, throttle, brake, and velocity respectively. The value of  $\lambda_s$ ,  $\lambda_t$ ,  $\lambda_b$ , and  $\lambda_v$  we used are listed in Table 9. For reinforcement learning, we adopt the training pipeline in DI-drive [2]<sup>1</sup>.

Hyper-parameter	Value
Initial learning rate	0.03
Weight decay	0.0001
Train epoch number	100
Train batch size	256
Key dictionary size	40960
$\lambda_{ins}$	1
$\lambda_{act}$	1
$\alpha$ for momentum update.	0.999
Hidden size of first layer in Projector	128
Hidden size of second layer in Projector	128
$\epsilon$ for action threshold	0.05

Table 7. ACO

Table 8. Inverse dynamics model

Hyper-parameter	Value
Learning rate	0.0003
Weight decay	0.0001
Train epoch number	50
Train batch size	256

### F Sampled frames of YouTube Driving Datasets

We provide sampled frames of five cities in Figure 9, 10, 11, 12, 13. Note that in captions we describe corresponding weather with blue text, region with red text, and scene style with green text. The columns are organized by different predicted steering values.

<sup>&</sup>lt;sup>1</sup> DI-drive [2] is an open-source auto-driving platform, supporting reinforcement learning in Carla. The reward scheme can be found at: https://github.com/opendilab/ DI-drive/blob/main/core/envs/simple\_carla\_env.py.



Fig. 9. Early morning drive at the center of Amsterdam, empty traffic on straight roads, with urban buildings surrounded



Fig. 10. Afternoon drive along the Italian Riviera, along narrow roads with frequent turnings



Fig. 11. Afternoon drive around the traditional Asakusa neighborhood, through rich traditional streets. Left-hand riding with middle density traffic



Fig. 12. Night drive in Madrid Spain, through Straight roads with middle density traffic



Fig. 13. Sunset drive in Lan Kwai Fong neighborhood, Hong Kong, middle traffic, modern building, left-hand, some one-way streets

Hyper-parameter	Value
Learning rate	0.0001
Weight decay	0.0001
Train epoch number	100
Train batch size	128
Velocity loss weight $\lambda_v$	0.05
Steering loss weight $\lambda_s$	0.5
Throttle loss weight $\lambda_t$	0.45
Brake loss weight $\lambda_b$	0.05

 Table 9. Imitation learning (CILRS [1])

 Table 10. Reinforcement learning (PPO [5])

Hyper-parameter	Value
Learning rate	0.0003
Weight decay	0.00001
SGD epoch number	6
$\lambda$ for GAE [4]	0.95
Train batch size	1024
SGD mini batch size	256
Clip ratio	0.2

# References

- Codevilla, F., Santana, E., López, A.M., Gaidon, A.: Exploring the limitations of behavior cloning for autonomous driving. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9329–9338 (2019)
- 2. drive Contributors, D.: DI-drive: OpenDILab decision intelligence platform for autonomous driving simulation. https://github.com/opendilab/DI-drive (2021)
- He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9729–9738 (2020)
- Schulman, J., Moritz, P., Levine, S., Jordan, M., Abbeel, P.: High-dimensional continuous control using generalized advantage estimation. arXiv preprint arXiv:1506.02438 (2015)
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347 (2017)