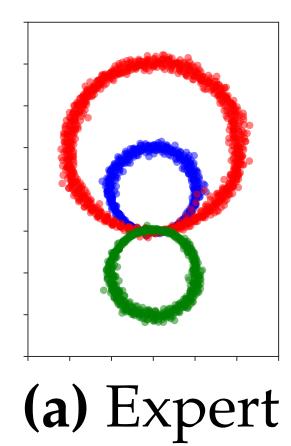
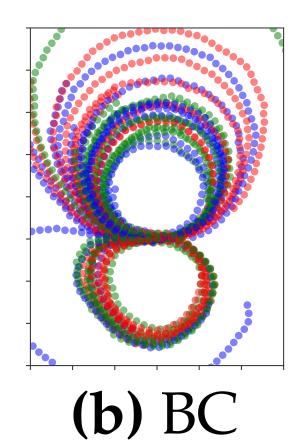
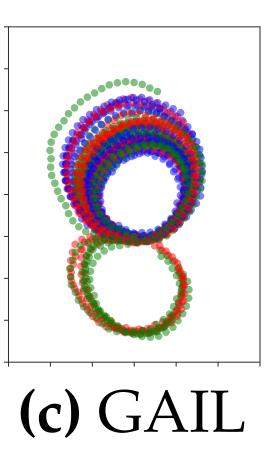


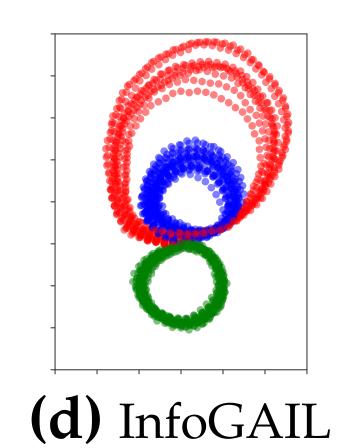
Introduction

- Imitation Learning mimic expert behavior without access to an explicit reward signal.
- Expert demonstrations provided by humans, however, often show significant variability.



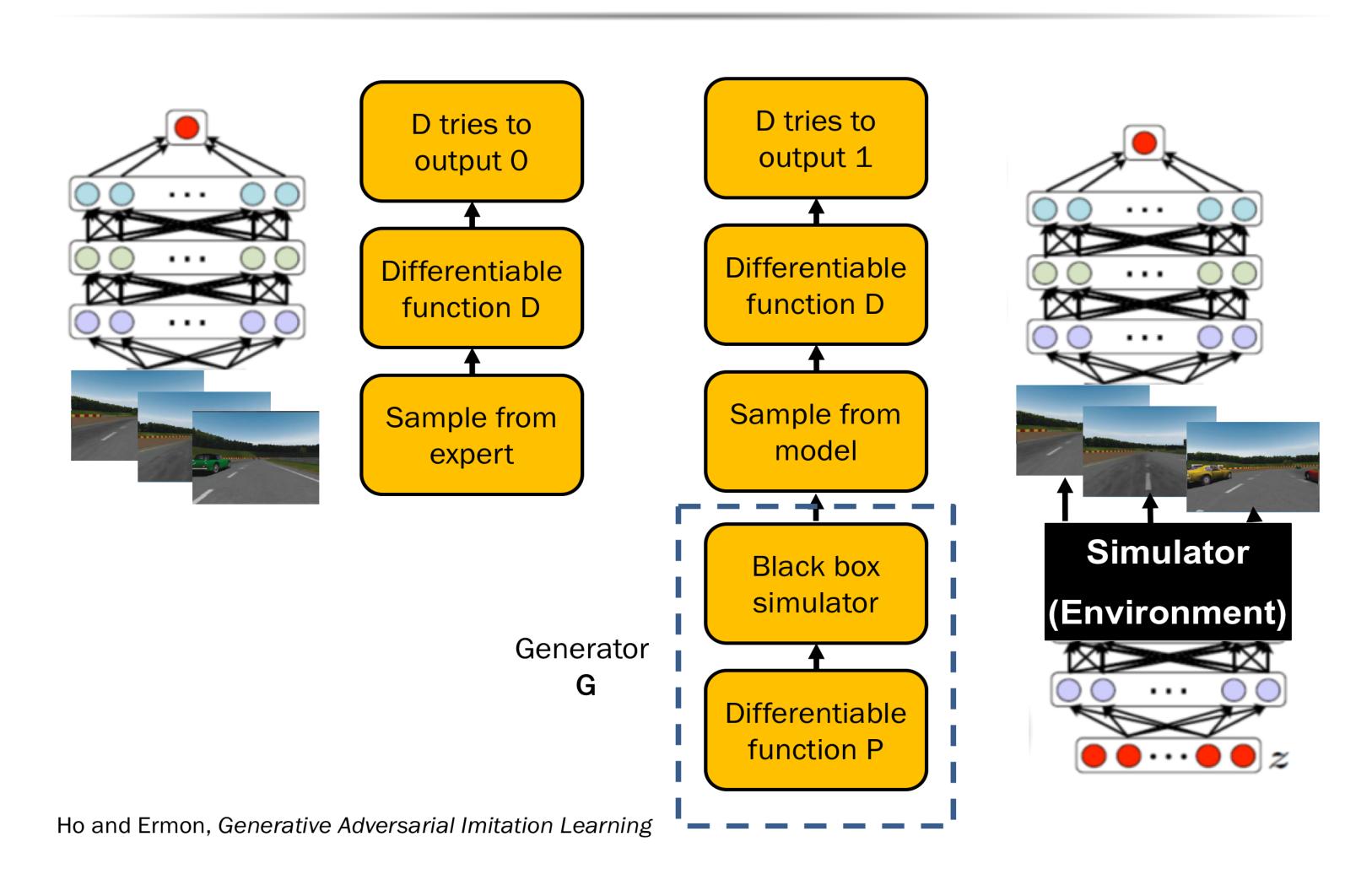






- BC deviates due to compounding errors.
- GAIL fails to capture the latent structure.
- Our method
- Can disentangle different behaviors (modes).
- Can do imitation learning from raw images.
- Can be used to anticipate actions.

Generative Adversarial Imitation Learning



- Discriminator *D* tries to distinguish between expert trajectories and ones from policy π .
- Policy π tries to fool the discriminator. $\min_{\pi} \max_{D} \mathbb{E}_{\pi}[\log D(s, a)] + \mathbb{E}_{\pi_{E}}[\log(1 - D(s, a))] \quad (1)$

InfoGAIL: Interpretable Imitation Learning from Visual Demonstrations

Yunzhu Li¹, Jiaming Song² and Stefano Ermon²

¹MIT Computer Science and Artificial Intelligence Laboratory ²Stanford Artificial Intelligence Laboratory

(2)

Interpretable Imitation Learning

Introduce a latent variable/code c.

• $L_{I}(\pi, Q)$ is a variational lower bound of MI $L_{I}(\pi, Q) = \mathbb{E}_{c \sim p(c), a \sim \pi(\cdot|s,c)}[\log Q(c|\tau)] + H(c)$ (3) where Q is an approximation to the posterior.

Algorithm 1 InfoGAIL

Input: Initial parameters of policy, discriminator and posterior approximation $\theta_0, \omega_0, \psi_0$; expert trajectories $\tau_E \sim \pi_E$ containing state-action pairs.

Output: Learned policy π_{θ}

for i = 0, 1, 2, ... do Sample latent codes: $c_i \sim p(c)$ Sample trajectories: $\tau_i \sim \pi_{\theta_i}(c_i)$. Sample state-action pairs $\chi_i \sim \tau_i$ and $\chi_E \sim \tau_E$.

Update ω_i to ω_{i+1} by ascending with gradients

 $\Delta_{\omega_i} = \hat{\mathbb{E}}_{\chi_i} [\nabla_{\omega_i} \log D_{\omega_i}(s, a)] + \hat{\mathbb{E}}_{\chi_E} [\nabla_{\omega_i} \log(1 - D_{\omega_i}(s, a))]$ Update ψ_i to ψ_{i+1} by descending with gradients

 $\Delta_{\psi_i} = -\lambda_1 \hat{\mathbb{E}}_{\chi_i} [\nabla_{\psi_i} \log Q_{\psi_i}(c|s,a)]$

Update θ_i to θ_{i+1} , using TRPO with the objective:

 $\mathbb{\hat{E}}_{\chi_{i}}[\log D_{\omega_{i+1}}(s,a)] - \lambda_{1}L_{I}(\pi_{\theta_{i}}, Q_{\psi_{i+1}})$

end for

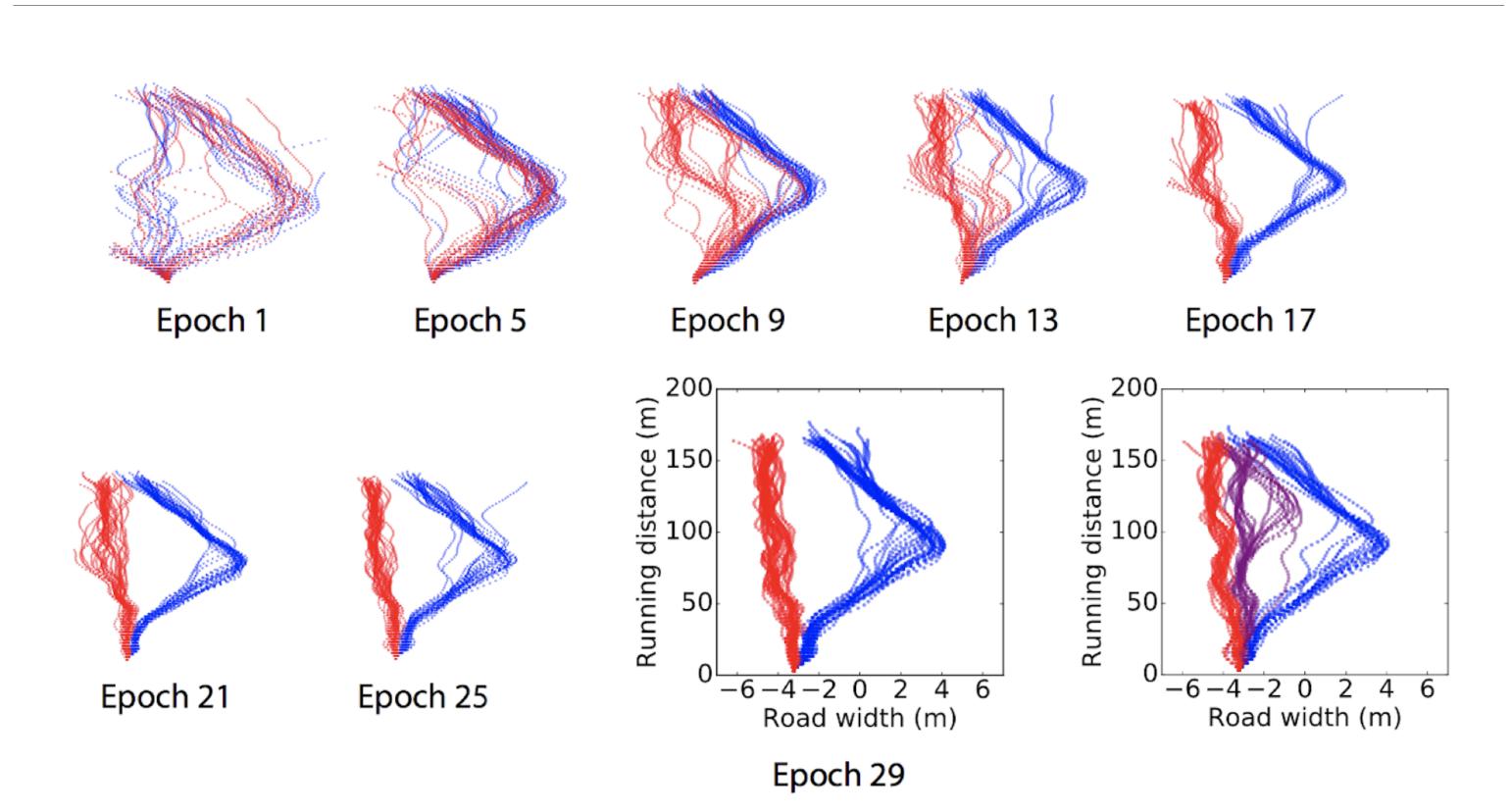


Figure: Visualization of different training stages.

InfoGAIL Training

- Reward Augmentation Incorporate prior knowledge by adding state-based incentives $\eta(\pi_{\theta}) = \mathbb{E}_{s \sim \pi_{\theta}}[r(s)].$
- Improved Objective
 Using WGAN to alleviate the problems of
 (1) vanishing gradient
 (2) mode collapse
 min max E_{πθ}[D_ω(s, a)] E_{πE}[D_ω(s, a)] λ₀η(π_θ)
 -λ₁L_I(π_θ, Q_ψ)
- Variance Reduction
 Baselines, Replay Buffers, etc.
- Transfer Learning for Visual Inputs
 Extract features from a pre-trained network.

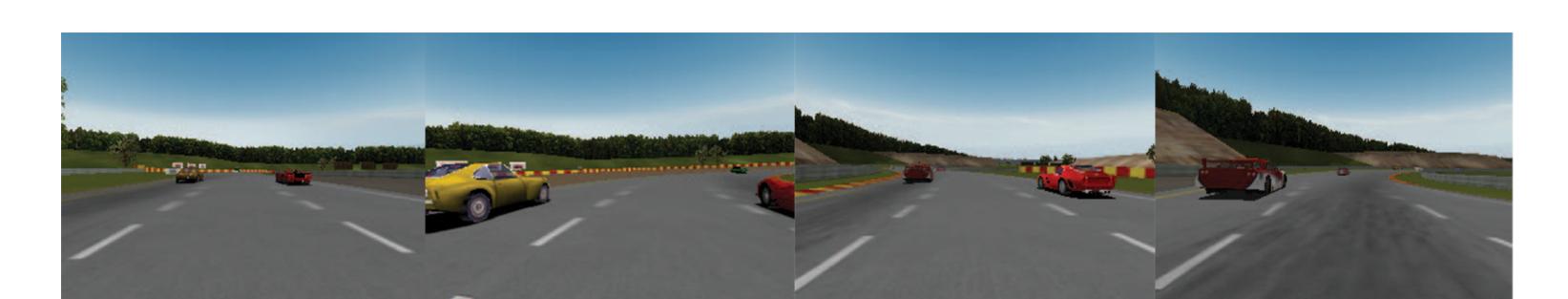


Figure: Transfer learning for handling visual inputs.

Experiments on Self-Driving

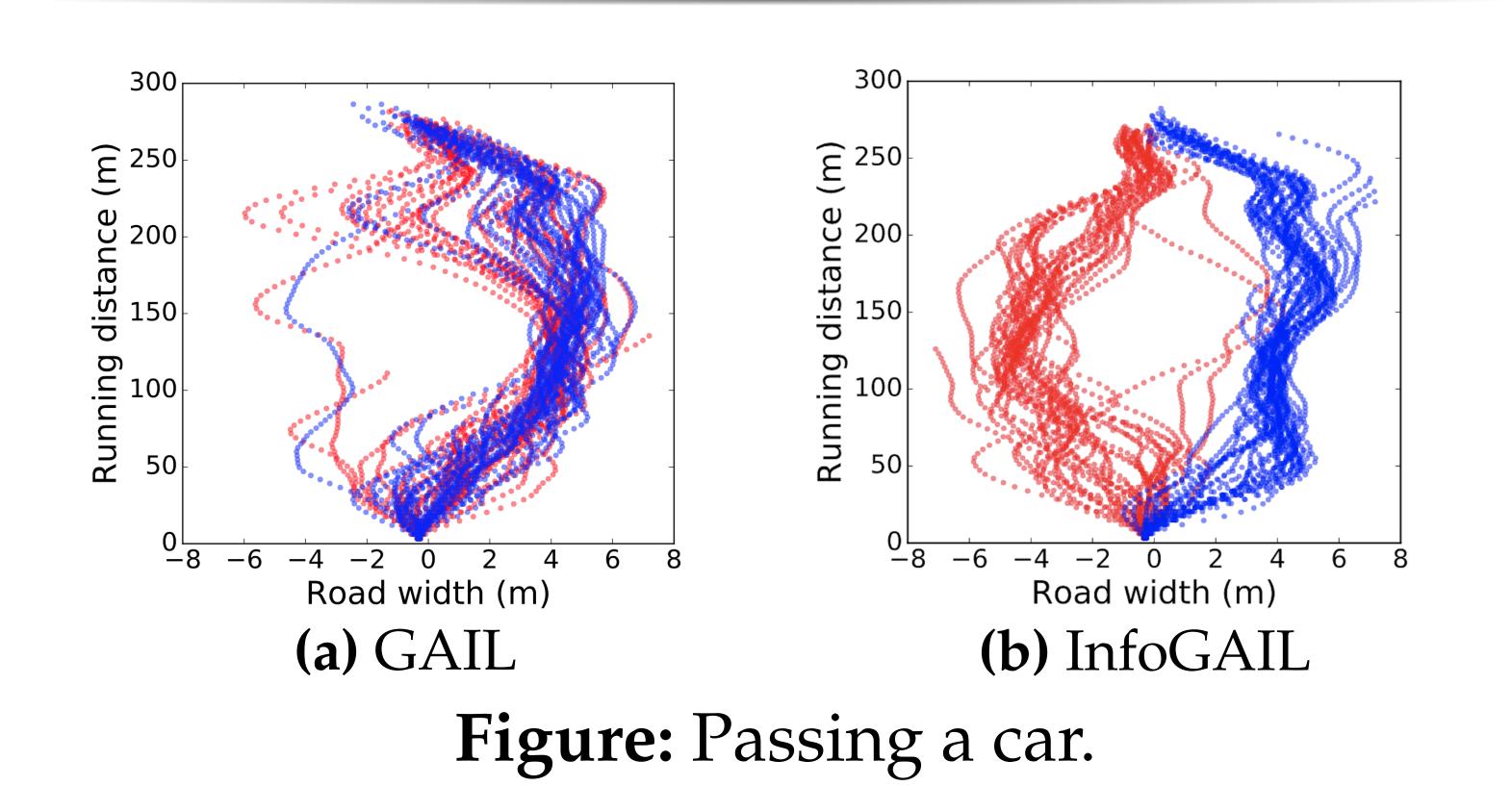
Interpretable Imitation Learning via Vision

- Using TORCS a driving simulator
- Vision as only source of perceptual inputs

The learned policy

- successfully distinguishes expert behaviors.
- produces interpretable representations from high-dimensional visual behavioral data.
- imitates each mode accordingly.
- low-level actions controlled by specifying high-level latent codes.





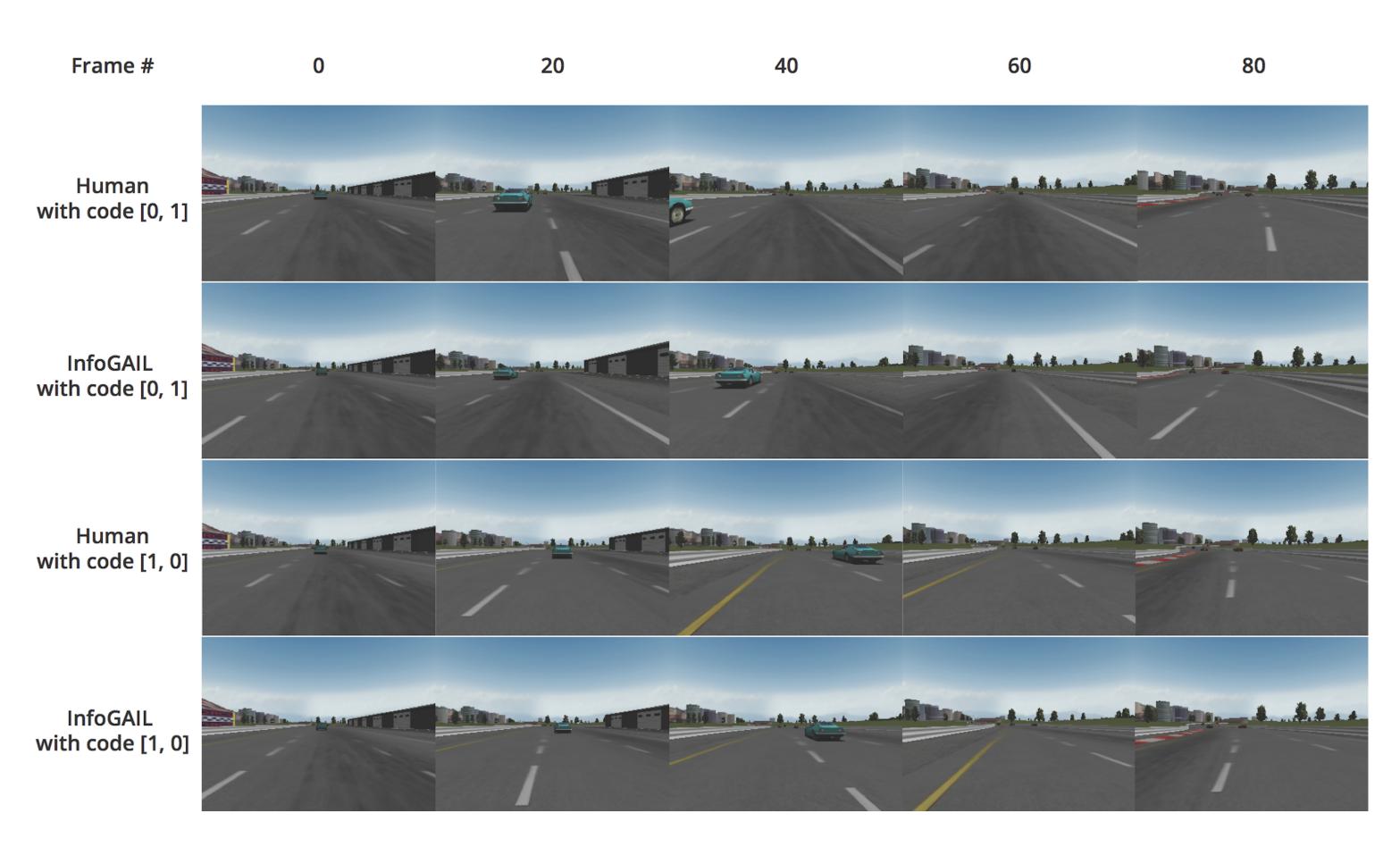


Figure: Visual inputs used for passing a car.

Table: Predictive

accuracy

(4)

Method	Acc.
Chance	50%
K-means	55.4%
PCA	61.7%
InfoGAIL (Ours)	81.9%
SVM	85.8%
CNN	90.8%

Table: Ablation study

Rollout dist.
701.83
914.45
1031.13
1123.89
1177.72
1226.68
1203.51

Code: https://github.com/ermongroup/infogail

Acknowledgements

Toyota Research Institute (TRI) provided funds to assist the authors with their research. This research was also supported by Intel Corporation, FLI and NSF grants 1651565, 1522054, 1733686.