

Learning Robust Driving Policies

Andreas Geiger

Autonomous Vision Group
University of Tübingen / MPI for Intelligent Systems Tübingen

August 23, 2020



University of Tübingen
MPI for Intelligent Systems

Autonomous Vision Group



Collaborators



Eshed Ohn-Bar



Aditya Prakash



Kashyap Chitta

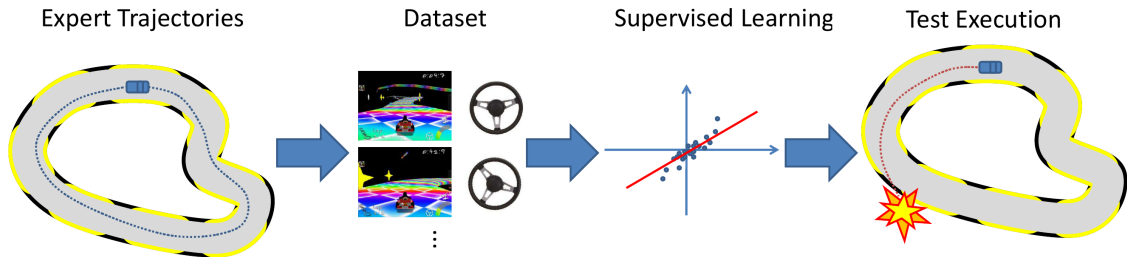


Aseem Behl



Andreas Geiger

Imitation Learning



Motivation: Hard coding policies is difficult \Rightarrow follow data-driven approach!

- ▶ **Given:** demonstrations or demonstrator
- ▶ **Goal:** train a policy to mimic decision

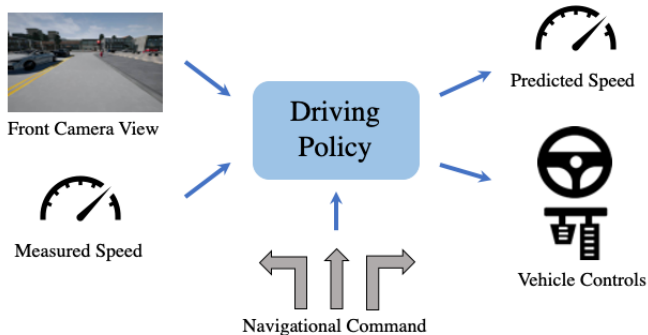
Conditional Imitation Learning

Advantages:

- ▶ End-to-End Trainable
- ▶ Cheap Annotations

Limitations:

- ▶ Generalization
- ▶ High Sample Complexity
- ▶ Covariate Shift
- ▶ Interpretability

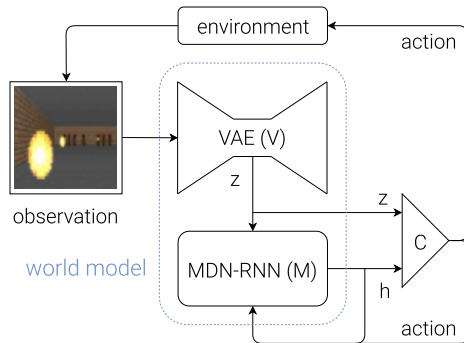


How can we learn to drive under the **vast diversity**
of all visual, planning and control scenarios?

Situational Driving

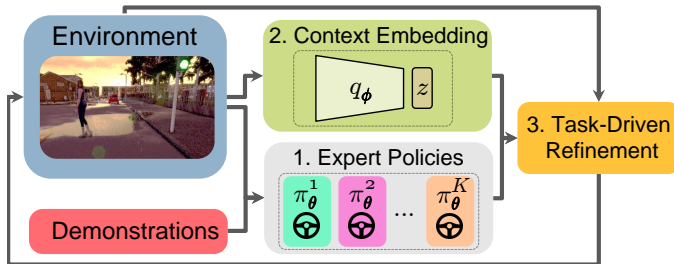


Inspiration: World Models



- Step 1: Learn **generative model** of game environments (VAE)
- Step 2: Learn dynamics model and control model in **latent space** (CMA-ES)
- Not sufficient \Rightarrow we combine this idea with **imitation learning**

Learning Situational Driving



- Step 1: Learn a **mixture of expert policies** $\{\alpha_\theta^k, \pi_\theta^k\}$ via imitation (LSD)
- Step 2: Learn a **general purpose context embedding** q_ϕ as a β -VAE
- Step 3: Perform **task-driven policy refinement** by interacting with the simulation and maximizing a driving task reward (LSD+)

Learning Situational Driving

$$\pi_{\Theta}(\mathbf{a}|\mathbf{o}, c) = \sum_{k=1}^K \underbrace{\alpha_{\theta}^k(\mathbf{o}, c)}_{\text{Mixture Weights}} \underbrace{\pi_{\theta}^k(\mathbf{a}|\mathbf{o}, c)}_{\text{Expert Models}} + \Psi \underbrace{\begin{bmatrix} q_{\phi}(\mathbf{I}) \\ v \\ c \end{bmatrix}}_{\text{Context Embedding}}$$
$$\pi_{\theta}^k(\mathbf{a}|\mathbf{o}, c) = \mathcal{N}\left(\mathbf{a} \mid \boldsymbol{\mu}_{\theta}^k(\mathbf{o}, c), \text{diag}(\boldsymbol{\sigma}_{\theta}^k(\mathbf{o}, c))^2\right)$$

Observations: $\mathbf{o} = [\mathbf{I}, v] \in \mathcal{O}$

Command: $c \in \mathcal{C} = \{\text{left, right, straight, follow}\}$

Actions: $\mathbf{a} \in \mathcal{A} = [-1, 1]^2$

Learning Situational Driving

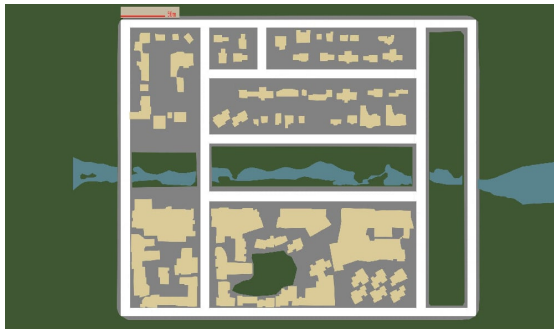
$$\pi_{\Theta}(\mathbf{a}|\mathbf{o}, c) = \sum_{k=1}^K \underbrace{\alpha_{\theta}^k(\mathbf{o}, c)}_{\text{Mixture Weights}} \underbrace{\pi_{\theta}^k(\mathbf{a}|\mathbf{o}, c)}_{\text{Expert Models}} + \underbrace{\Psi \begin{bmatrix} q_{\phi}(\mathbf{I}) \\ v \\ c \end{bmatrix}}_{\text{Context Embedding}}$$
$$\pi_{\theta}^k(\mathbf{a}|\mathbf{o}, c) = \mathcal{N} \left(\mathbf{a} \mid \boldsymbol{\mu}_{\theta}^k(\mathbf{o}, c), \text{diag}(\boldsymbol{\sigma}_{\theta}^k(\mathbf{o}, c))^2 \right)$$

Training:

- Step 1: Learn Mixture of Experts: $\mathcal{L}_{\text{MoE}} = -\log \left[\sum_{k=1}^K \alpha_{\theta}^k \pi_{\theta}^k \right] + \mathcal{L}_V + \mathcal{L}_R$
- Step 2: Learn Context Embedding: $\mathcal{L}_{\text{VAE}} = \beta \text{KL} (q_{\phi}(\mathbf{z}|\mathbf{I}) \parallel p_0(\mathbf{z})) + \|d_{\phi}(\mathbf{z}) - \mathbf{I}\|_2^2$
- Step 3: Task-driven optimization: $\mathcal{J}_{\text{TASK}}(\boldsymbol{\theta}_{\text{readout}}, \boldsymbol{\Psi}) = \mathbb{E}_{\pi_{\Theta}} \left[\sum_{t=0}^T r_t \right]$

Experiments

CARLA Benchmark



Training Town



Test Town

- ▶ Random start and end location, 4 known weathers, 2 unseen weathers
- ▶ Metric: Percentage of successfully completed episodes (success rate)
- ▶ Collision does not necessarily terminate episode

NoCrash Benchmark



Empty



Regular



Dense

- ▶ Difficulty varies with number of dynamic agents in the scene
- ▶ Empty: 0 Agents Regular: 65 Agents Dense: 220 Agents
- ▶ All collisions terminate episode

AnyWeather Benchmark



- Evaluation on 10 unseen weathers, quantifies generalization performance

Importance of Mixture Model

Evaluation Task	Training Data and Mixture Components		
	Navigation (Static, K=1)	Navigation (Dynamic, K=1)	Navigation (Dynamic, K=3)
Straight (Static)	99	64	100
One Turn (Static)	98	74	100
Navigation (Static)	96	78	98
Navigation (Dynamic)	40	78	92

Results of Mixture Model on CARLA Benchmark:

- Static model solves static scenes well but cannot handle dynamic objects
- Dynamic model handles dynamic scenes better but degrades on static scenes
- Dynamic mixture model generalizes to all scenarios (without on-policy data)

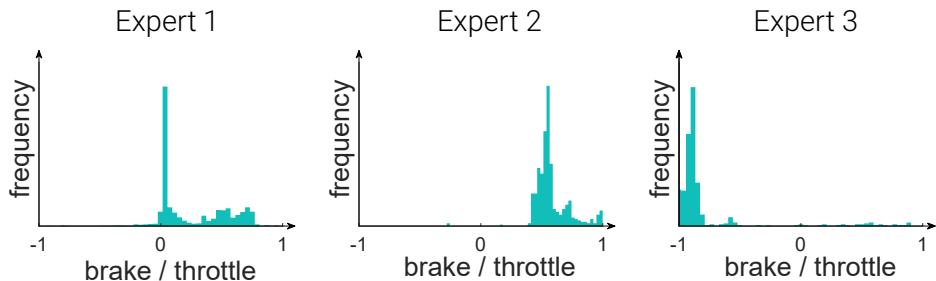
Importance of Mixture Model and Task-based Refinement

Model	Success Rate (%)
Monolithic (K=1)	75
MoE Shared Backbone (K=3)	89
MoE Shared Backbone (K=5)	90
MoE Shared Backbone (K=8)	87
MoE Separate Backbone (K=3)	94
MoE Separate Backbone (K=5)	93
MoE Separate Backbone (K=8)	93
MoE Separate Backbone + Refinement (K=3)	98

Results of Full Model on CARLA Benchmark:

- Performance improves up to 3 or 5 mixture components
- Separate backbones increase diversity and generalization
- Tasked-based refinement improves performance further

Emergent Driving Modes



Emergent Driving Modes:

- Acceleration distribution of three different experts during testing

Results on CARLA Benchmark

Driving Task	CIRL	CILRS	CILRS (ours)	LSD (ours)	LSD+R (ours)
Straight	100	96	96	100	100
One Turn	71	84	86	99	99
Navigation	53	69	67	99	99
Navigation Dynamic	41	66	64	94	98

- Using reward-based optimization alone (CIRL) is not sufficient
- LSD enables better driving behavior across all driving tasks
- Large improvements in the presence of dynamic objects

Results on CARLA NoCrash Benchmark

Driving Task	CILRS	CILRS	LSD (ours)	LSD+R (ours)	Expert
Empty	66 ± 2	65 ± 2	93 ± 2	94 ± 1	96 ± 0
Regular	49 ± 5	46 ± 2	66 ± 2	68 ± 2	91 ± 1
Dense	23 ± 1	20 ± 1	27 ± 2	30 ± 4	41 ± 2

- ▶ All methods perform worse due to challenges (density, collision terminations)
- ▶ Expert provided by CARLA often fails in dense environments (e.g., clogging)
- ▶ LSD enables better driving behavior across all driving tasks

Results on AnyWeather Benchmark

Task	CILRS	LSD (ours)	LSD+R (ours)
Straight	83.2	85.2	85.6
One Turn	78.4	80.4	81.6
Navigation	76.4	78.8	79.6
Nav. Dynamic	75.6	77.2	78.4

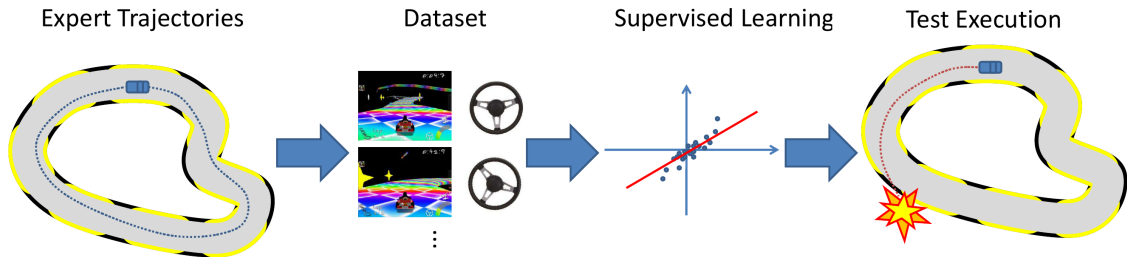
- ▶ AnyWeather benchmark test generalization to challenging unseen weathers
- ▶ All methods can fail even on simple straight driving tasks
- ▶ Some challenging weathers lead to zero success rate for all methods
- ▶ More research is required to address these challenges

Qualitative Results



How useful is **data aggregation** for self-driving?

Imitation Learning



Hard coding policies is often difficult \Rightarrow Rather use a data-driven approach!

- **Given:** demonstrations or demonstrator
- **Goal:** train a policy to mimic decision

Formal Definition of Imitation Learning

General Imitation Learning:

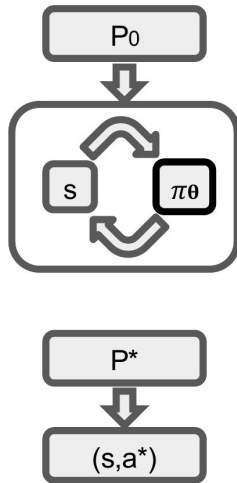
$$\operatorname{argmin}_{\theta} \mathbb{E}_{s \sim P(s|\pi_{\theta})} [\mathcal{L}(\pi^*(s), \pi_{\theta}(s))]$$

- State distribution $P(s|\pi_{\theta})$ depends on rollout determined by current policy π_{θ}

Behavior Cloning:

$$\operatorname{argmin}_{\theta} \underbrace{\mathbb{E}_{(s^*, a^*) \sim P^*} [\mathcal{L}(a^*, \pi_{\theta}(s^*))]}_{= \sum_{i=1}^N \mathcal{L}(a_i^*, \pi_{\theta}(s_i^*))}$$

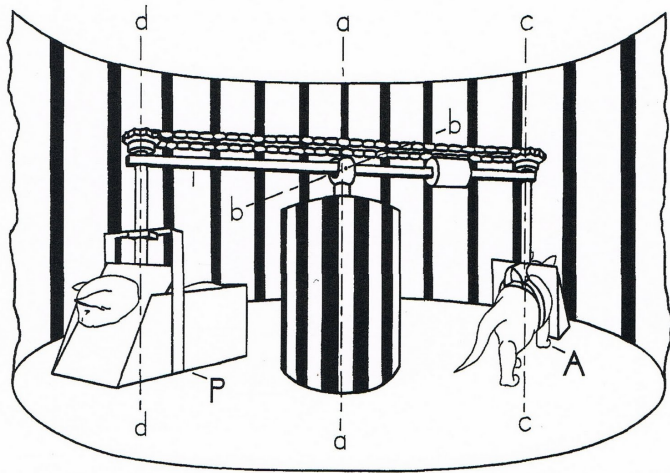
- State distribution P^* provided by expert
- Reduces to supervised learning problem



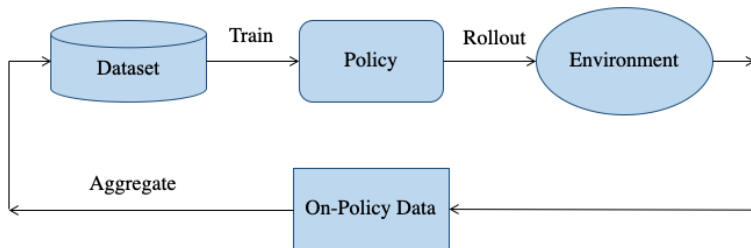
Challenges of Behavior Cloning

- ▶ Behavior cloning makes IID assumption
 - ▶ Next state is sampled from states observed during expert demonstration
 - ▶ Thus, next state is sampled independently from action predicted by current policy
- ▶ What if π_θ makes a mistake?
 - ▶ Enters new states that haven't been observed before
 - ▶ New states not sampled from same (expert) distribution anymore
 - ▶ Cannot recover, can lead to catastrophic failure

Experiment by Held and Hein



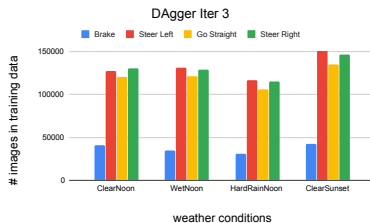
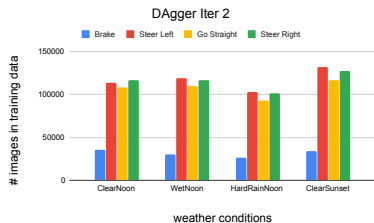
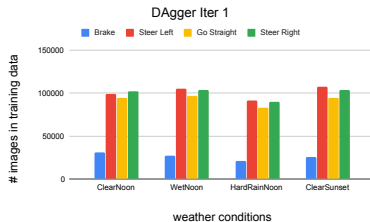
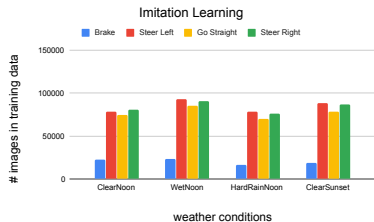
DAgger



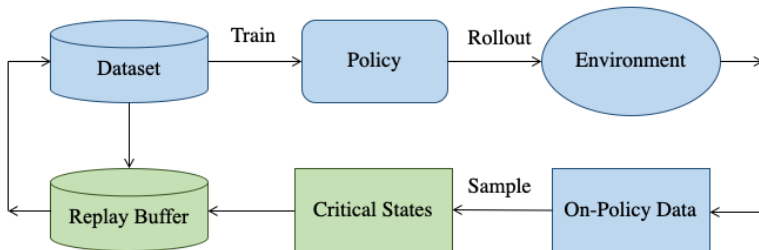
Data Aggregation (DAgger):

- ▶ Iteratively build a set of inputs that the final policy is likely to encounter based on previous experience. Query expert for aggregate dataset.
- ▶ But can easily overfit to main mode of demonstrations
- ▶ High training variance (random initialization, order of data)

Distribution over Driving Actions



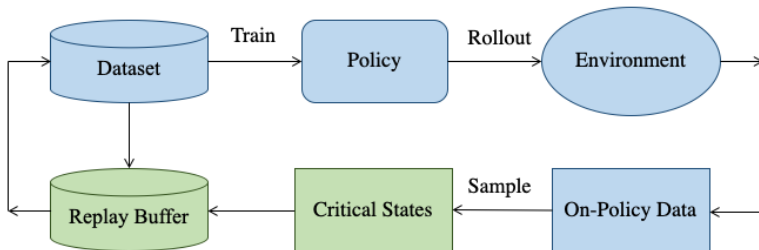
Dagger with Critical States and Replay Buffer



Key Ideas:

1. Sample **critical states** from the collected on-policy data based on the utility they provide to the learned policy in terms of driving behavior
2. Incorporate a **replay buffer** which progressively focuses on the high uncertainty regions of the policy's state distribution

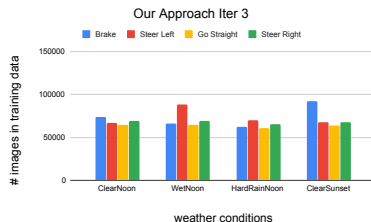
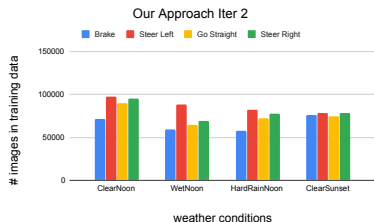
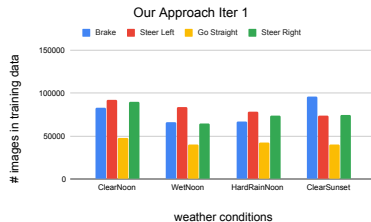
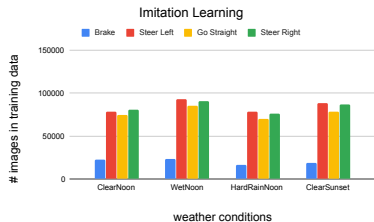
Dagger with Critical States and Replay Buffer



Sampling Strategies:

- ▶ Task-based: Sample uniformly from “left”, “right”, “straight”
- ▶ Policy-based: Use test-time dropout to estimate epistemic uncertainty
- ▶ Expert-based: Highest loss or deviation in brake signal wrt. expert

Distribution over Driving Actions



Experiments

Evaluation



Empty



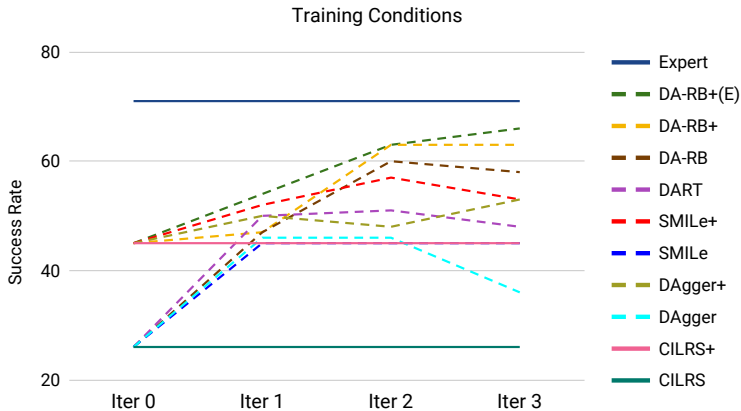
Regular



Dense

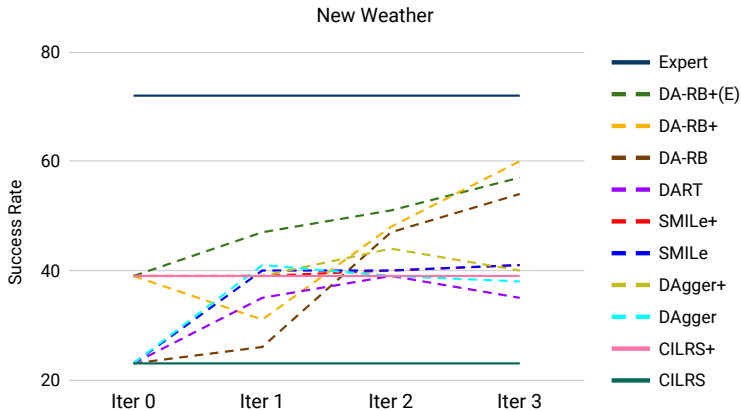
- ▶ CARLA **NoCrash benchmark**
- ▶ **Dense setting** with 220 agents
- ▶ Comparison to various baselines with (+) and without data augmentation

Evaluation



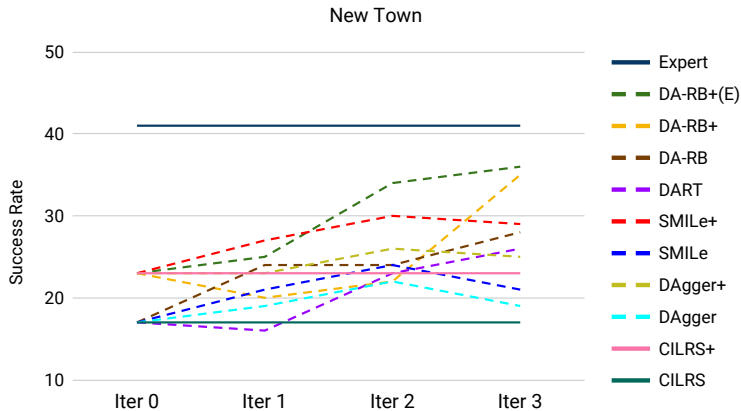
- Data augmentation increases the performance of all methods
- DAgger overfits quickly (!), not better than data augmentation
- Our model consistently improves upon the baselines in all conditions

Evaluation



- Data augmentation increases the performance of all methods
- DAgger overfits quickly (!), not better than data augmentation
- Our model consistently improves upon the baselines in all conditions

Evaluation



- Data augmentation increases the performance of all methods
- DAGger overfits quickly (!), not better than data augmentation
- Our model consistently improves upon the baselines in all conditions

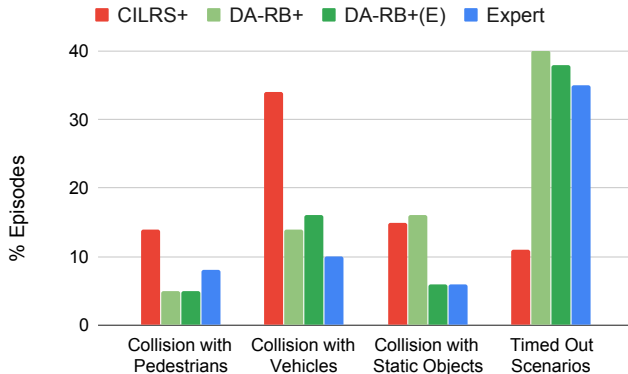
Evaluation

Task	CILRS ⁺	DART	DA-RB ⁺ (Ours)	DA-RB ⁺ (E) (Ours)	Expert
Training	45±6	50±1	62±1	66 ±5	71±4
New Weather	39±4	37±2	60 ±1	56±1	72±3
New Town	23±1	26±2	34±2	36 ±3	41±2
New Town & Weather	26±2	21±1	25±1	35 ±2	43±2

Mean and standard deviation over 3 evaluation runs.

- Ensemble (E) improves performance further (particularly in new environment)
- Expert provided by CARLA often fails in dense environments (e.g., clogging)

Infractions Analysis



- ▶ Significant reduction in collisions with dynamic objects
- ▶ More time-outs due to less infractions (e.g., clogged scenes, red lights)

Training Variance

	CILRS ⁺	DAgger ⁺	DA-RB ⁺
Iter 0	14.6 \pm 3.4	14.6 \pm 3.4	14.6 \pm 3.4
Iter 1	-	15.2 \pm 5.1	24.8 \pm 1.9
Iter 2	-	13.2 \pm 1.9	25.4 \pm 1.5
Iter 3	-	17.8 \pm 3.6	27.0 \pm 0.9

Standard deviation wrt. 5 random training seeds (New Town & Weather)

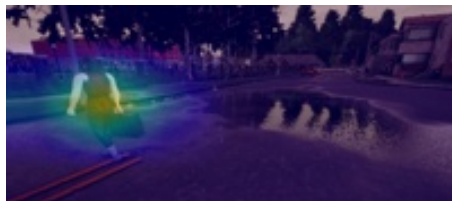
- Significant reduction in variance compared to CILRS and DAgger
- Sampling the dataset based on critical states is crucial for variance reduction

Interpretability: GradCAM Attention Maps

CILRS [Codevilla et al. 2019]



Our Approach



Interpretability: GradCAM Attention Maps

CILRS [Codevilla et al. 2019]



Our Approach



Qualitative Results

CILRS+ (Codevilla et al. 2019)



DA-RB+ (Our Approach)



Summary

Summary

- ▶ A single imitation learner cannot capture the complexities of driving
- ▶ Mixture of experts can significantly improve generalization
- ▶ Task-driven optimization is difficult but important
- ▶ Data augmentation is important but can easily overfit in self-driving
- ▶ Sampling critical states in a replay buffer improves aggregation performance
- ▶ Training variance can be reduced using this strategy
- ▶ Generalization to all CARLA weathers remains unsolved
- ▶ Better experts and improvements in the CARLA simulation are needed
(and on their way – current version is CARLA 0.9.9!)

Thank you!

<http://autonomousvision.github.io>



European Research Council
Established by the European Commission



Federal Ministry
of Education
and Research



Microsoft®
Research

