

Augmenting Imitation Experience via Equivariant Representations

Dhruv Sharma^{*†} Alihusein Kuwajerwala^{*†} Florian Shkurti[†]

Abstract—The robustness of visual navigation policies trained through imitation often hinges on the augmentation of the training image-action pairs. Traditionally, this has been done by collecting data from multiple cameras, by using standard data augmentations from computer vision, such as adding random noise to each image, or by synthesizing training images. In this paper we show that there is another practical alternative for data augmentation for visual navigation based on extrapolating viewpoint embeddings and actions nearby the ones observed in the training data. Our method makes use of the geometry of the visual navigation problem in 2D and 3D and relies on policies that are functions of *equivariant* embeddings, as opposed to images. Given an image-action pair from a training navigation dataset, our neural network model predicts the latent representations of images at nearby viewpoints, using the equivariance property, and augments the dataset. We then train a policy on the augmented dataset. Our simulation results indicate that policies trained in this way exhibit reduced cross-track error, and require fewer interventions compared to policies trained using standard augmentation methods. We also show similar results in autonomous visual navigation by a real ground robot along a path of over 500m.

I. INTRODUCTION

From autonomous driving to flying and manipulation, vision-based policies learned through imitation are starting to be deployed on robot systems of critical importance. To ensure safety during deployment, we would ideally like to have formal guarantees about the generalization and robustness properties of these policies. While desirable, however, formal guarantees for out-of-distribution generalization have been attained on perturbation models that are theoretically convenient, but limited in practice [1], as they often do not capture the full range of plausible unseen states. Data augmentation during training is significant in practice.

In this paper we address this problem by proposing a dataset augmentation scheme for training policies whose inputs are image embeddings, as shown in Fig. 1. Given an observed image from the center camera of a vehicle, our method allows us to learn a map that predicts the embeddings of corresponding nearby viewpoints, as well as the actions that would have needed to be taken if the vehicle’s camera was at those viewpoints. Our method relies on predicting embeddings of nearby viewpoints using *equivariant mappings* [2], [3] that are trained to transform the embedding of the center camera to the embeddings of nearby viewpoints, both for ground vehicles and for flying vehicles, as we will show in the evaluation section. Equivariance has emerged as

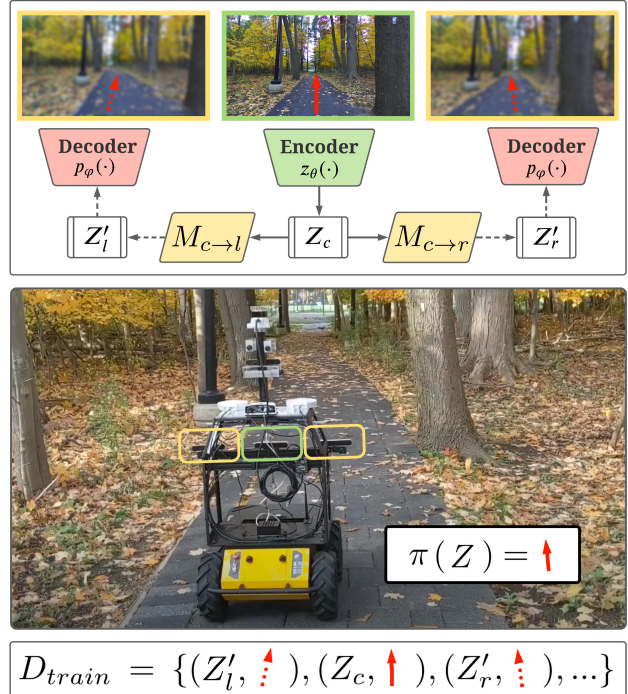


Fig. 1. An overview of our system. Images from the center camera are embedded into a low-dimensional representation, Z_c , which is deterministically transformed into the predicted embeddings of the left and right viewpoints respectively, via equivariant mappings M . Both the actual steering command (solid red arrow) and the predicted steering commands (dashed red arrows) are included in the augmented training set on which the policy is trained. At test time, the policy is executed on embeddings from the center camera. The maps M are trained separately from π .

an important structural bias for neural network design in the last few years. Here we use it as an auxiliary loss to the main behavioral cloning loss.

Our main contribution is showing that augmenting the training dataset with predicted embeddings of nearby viewpoints, using learned equivariant maps, increases the robustness of the vision-based policy. This results in lower cross track error and fewer human interventions needed for navigation tasks, both for ground vehicle and aerial vehicle navigation tasks in 2D and 3D respectively. We demonstrate that our method applies to real world visual navigation scenarios by deploying it on a terrestrial mobile robot, resulting in significantly improved navigation performance.

II. RELATED WORK

Equivariant representations: Equivariance for image representations captures relationships between encodings of two related images that are determined by the transformation of those images [3]. If the encodings are the same, then they

^{*}Authors contributed equally.

[†]Robot Vision and Learning Laboratory, University of Toronto Robotics Institute. {dhruv, florian}@cs.toronto.edu, ali.kuwajerwala@mail.utoronto.ca

are invariant to that transformation. In our work, we want the representation of the left camera to be predictable by the center camera, so we want a high degree of equivariance. Feature equivariance and invariance in terms of rotation transformations was studied in [4], [5], [6], as well as in the context of probabilistic models [7], [8], [9]. The utility of equivariant representations for multiple downstream tasks, learned by egomotion was recognized in [2]. Our work here is related the most to this paper.

Equivariant representations have also been useful for 3D volumetric data [10], even in the fully unsupervised setting [11], addressing equivariance and invariance to both rotation and translation for 3D data [12], [13]. Additionally, instead of imposing equivariance through a loss function, as in [2], many recent works build equivariance in the structure of the network [14], [15], [16], [17], [18], [19], for example in the multi-view setting [20], [21], [22], [23], [24]. Finally, equivariance has also been used to improve state representations for planning in MDPs [25].

Automated data augmentations: Typical dataset augmentations for supervised deep learning have included addition of noise, rotations, crops, blurring, and others. While this is usually a set of manually hand-crafted augmentation schemes, there has been increased interest in automatically computing data augmentations [26], [27], [28], [29] to increase the robustness of classifiers, regressors, and RL policies. Augmentations have also been used in monocular depth estimation [30].

Robustness in behavioral cloning: Behavioral cloning techniques that learn policies through pure supervised learning from a given dataset of state-action pairs usually suffer from covariate shift [31], [32]. Dean Pomerleau’s PhD thesis [33] and pioneering work on ALVINN [34] details a data augmentation scheme that injects synthetic images of the road ahead from viewpoints nearby the camera to the training set. Our paper provides an alternative that is based on neural network view synthesis, which in the last years has seen significant progress. Pomerleau’s thesis also identified issues of covariate shift, which methods such as DAgger [35] address by iteratively labeling the newly visited states after each policy evaluation. The question of dataset augmentation in one round of behavioral cloning, however, remains. Collecting large driving datasets, in a way that no manual annotation is required from the user, is critical for autonomous visual navigation.

One example of this was in UAV navigation in forest trails [36]. The images of the center camera are automatically annotated with the forward steering command, the images of the left camera with the right, and the images of the right camera are annotated with the left steering command. A similar hardware setup and automated data collection was used by the NVIDIA self-driving car team [37], as well as more recent work on behavior cloning conditioned on human commands (e.g. turn left at the next intersection) [38]. The crux behind the idea of the 3-camera setup is that behavioral cloning can become more robust if we obtain expert actions from states that are nearby the forward-facing

one. Our method builds upon this approach by replacing the need for additional sensors with a model that predicts their embeddings, whose corresponding actions are automatically annotated. We note that in other similar approaches, such as DART [39] an expert demonstrator will have to manually annotate the additional states. In DART, Gaussian noise is added to the original training dataset in order to explicitly include nearby labeled states in the augmented dataset, together with manual annotations.

III. METHODOLOGY

Our goal is to learn an encoding method that satisfies equivariance constraints with respect to camera poses or viewpoints nearby the original camera. This allows us to predict the embeddings of nearby viewpoints from the observations of the center camera. Here, nearby viewpoints refer to horizontal and vertical translation of the camera over short distances capturing the same scene. Our case involves translations of 0.25 to 0.5 meters.

Given an image I_i of a scene taken from viewpoint i , we compute its corresponding embedding in the latent space Z_i , and use an equivariant map $M_{i \rightarrow j}$ to transform it to the embedding corresponding to another desired nearby viewpoint j . This predicted embedding Z'_j can then be used as augmented training data for policies trained on such embeddings and actions.

We train an image encoder network on a dataset of image pairs, where each pair consists of images of the same scene from neighboring viewpoints. In our case, the images are obtained from synchronized videos recorded by a set of cameras mounted on a vehicle visually navigating through an environment. The relative poses of the cameras are assumed to be known and looking at the same scene. With each image pair we simultaneously also record the corresponding expert control command used to navigate the vehicle at that time.

We then use the learned embeddings to train policies using imitation learning.

A. Equivariant Feature Learning

Given a tuple of images $\mathbf{I} = \langle I_1, I_2, \dots, I_n \rangle$ that capture the scene at the same timestep, our model learns an encoder $z_\theta(\cdot) : I_i \rightarrow \mathbb{R}^D$, parametrized by θ , that maps an image to an embedding Z_i . We want the embeddings to exhibit equivariance, i.e. change predictably with respect to transformations applied to the viewpoint i of the input image I_i . Let g be a transformation applied to the input image I in the pixel space, and let M_g define a corresponding transformation in the latent space. Then:

$$z_\theta(g(I)) \approx M_g(z_\theta(I)) \quad (1)$$

For example, suppose the images I_c and I_l correspond to the center and left camera images in Fig. 1, respectively. If g is an image transformation such that $g(I_c) = I_l$, then $M_g = M_{c \rightarrow l}$ defines a transformation in the latent space such that $M_{c \rightarrow l}(Z_c) = Z'_l \approx Z_l$. In latent space, we want Z_c and Z_l to maintain a geometric relation directly corresponding to the one between images I_c and I_l .

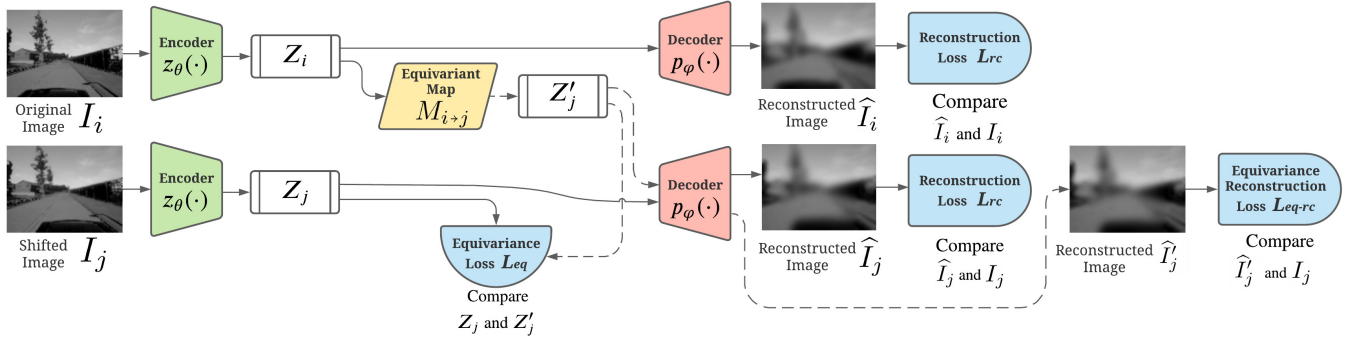


Fig. 2. Representation Learning Phase: we use an encoder-decoder network to learn image embeddings that are equivariant to input camera viewpoint. Mapping network $M_{i \rightarrow j}$ is learned to map these embeddings from viewpoint i to viewpoint j . The encoder-decoder and the mapping networks are trained collectively by optimizing a combined loss function shown in Eqn. 5. See III for details.

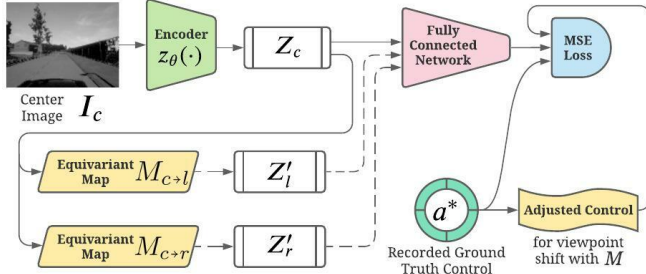


Fig. 3. Policy Training Phase: The center image is passed to an encoder, then the resulting embedding is mapped to nearby viewpoints (for experiments using an equivariant map), and finally the embeddings are passed to a FCN model to produce steering output. Note that the encoder and FCN models are trained *separately*. See III-D for details.

B. Equivariant Map Learning as an Auxiliary Task

Now, we design an objective function that encourages the learned embeddings to exhibit equivariance with respect to camera transformations, or camera viewpoints in our case. Given a tuple of images $\mathbf{I} = \langle I_1, I_2, \dots, I_n \rangle$ captured at the same time, where image I_i is taken from the i^{th} viewpoint, and the corresponding embeddings $z_\theta(I_i)$, we set up the following objective function to enforce equivariance between any two image viewpoints.

$$L_{eq}(\theta, M) = \sum_i \sum_j \|z_\theta(I_j) - M_{i \rightarrow j}(z_\theta(I_i))\|^2 \quad (2)$$

C. Image Reconstruction as an Auxiliary Task

We make use of image reconstruction as an auxiliary task, so that the embeddings capture the information needed to reconstruct images across all different viewpoints. We jointly train an encoder-decoder model for each viewpoint, as shown in Fig. 2. The encoder, decoder, as well as the equivariant mapping $M_{i \rightarrow j}$ are optimized together in a shared objective. The latent representations Z_i are decoded back to images by a decoder $p_\varphi(\cdot)$. The output of the decoder is then used to compute and minimize the autoencoding loss L_{rc} from the same viewpoint:

$$L_{rc}(\theta, \varphi) = \sum_i \|I_i - p_\varphi(z_\theta(I_i))\|^2 \quad (3)$$

To further reinforce the effectiveness of the equivariant maps $M_{i \rightarrow j}$, we enforce an additional reconstruction penalty

between pairs of viewpoints:

$$L_{eq-rc}(\theta, \varphi, M) = \sum_i \sum_j \|I_j - p_\varphi(M_{i \rightarrow j}(z_\theta(I_i)))\|^2 \quad (4)$$

The overall representation learning loss is a combination of the reconstruction, equivariant reconstruction, as well as the equivariant loss:

$$L(\theta, \varphi, M) = \lambda_1 L_{rc}(\theta, \varphi) + \lambda_2 L_{eq}(\theta, M) + \lambda_3 L_{eq-rc}(\theta, \varphi, M) \quad (5)$$

We backprop through the combined loss in Eqn. 5 to optimize θ , φ , and M . The hyperparameters are chosen as $\lambda_1 = 1$, $\lambda_2 = 10$, and $\lambda_3 = 1$ to ensure all loss terms have relatively same scale.

D. Imitation Learning with learned Features

After we learn $z_\theta(\cdot)$ using the representation learning loss in Eqn. 5, we fix the encoder weights and learn a control policy, i.e. a multi-layer perceptron representing $\pi(Z)$, via imitation learning. The control policy $\pi(\cdot)$ is trained on embedding-action pairs both from the center camera (observed embeddings and actions) and from the nearby viewpoints (predicted embeddings from equivariance relations and predicted actions). The objective function used is the following:

$$L(\pi) = \sum_i \|a_i - \pi(Z_i)\|^2 \quad (6)$$

At test time and deployment, the encoder is used to convert the images to embeddings which are then used as input to the learned policy, as shown in Fig. 2.

IV. EXPERIMENTS

We validate our approach in two different experimental settings, an aerial vehicle and a terrestrial vehicle. The experiments involve the following steps:

Data Collection. Each data sample $\langle \mathbf{I}, \mathbf{A} \rangle$ consists of an image tuple \mathbf{I} that consists of images from different camera viewpoints and an action set \mathbf{A} that contains the corresponding control actions. We collect data both in simulation and in the real world, for the respective experiments.

Learning Equivariant Features. We learn the feature mapping function $z_\theta(\cdot)$ and equivariant mappings $M_{i \rightarrow j}$

using the collected image sets I in the original training dataset.

Policy Evaluation. In order to assess the effectiveness of our learned embeddings and equivariant mappings we evaluate and compare three imitation learning policies that each learn to map images to control actions:

- 1) The first policy is trained on image-action pairs obtained from a single camera, in our case the center camera mounted on the vehicle.
- 2) The second policy is trained on image-action pairs obtained from all the camera sensors mounted on the vehicle. In our case it is 3 viewpoints for the car and 5 for the flying vehicle. The recorded expert control is modified to account for the shifted camera viewpoints for the off center cameras.
- 3) The third policy uses our equivariant augmentations and is trained on image-action pairs obtained from the center camera and predicted nearby embeddings. The center camera embedding is mapped to the embeddings of the other viewpoints in the latent space. This corresponds to the training scheme described in Sec. III-C.

Evaluation Metrics. To compare the three policies, we record the number of interventions (which occur when the ground vehicle goes off of the current lane or when the flying vehicle deviates more than 2.5 meters away from the goal trajectory), average cross track error (deviation from the specified trajectory in meters), and autonomy. Autonomy here is calculated as in [37], i.e. the percentage of time the car was not being driven by the learned policy due to an intervention. More specifically, for an experiment with n interventions:

$$\text{autonomy} = \left(1 - \frac{n \cdot (6 \text{ seconds})}{\text{elapsed time in seconds}}\right) \cdot 100 \quad (7)$$

The experiment setup details and results are discussed in Sec. IV-A and Sec. IV-B.

Policy and Encoder Architectures. The network architecture used for vehicle navigation was similar to the one used by NVIDIA in [37], it consists of a CNN encoder network followed by a FCN regressor network. An important distinction between their work and ours is that the entire model in [37] is trained end-to-end, whereas we train the encoder separately from the MLP, as mentioned in Sec III-D. For consistent comparisons, we ensure that the encoder $z_\theta(I)$ has the same architecture as the image encoder in [37].

Policy Training. All the networks are optimized using ADAM optimizer [40] with a learning rate of 10^{-4} and a batch size of 64. We use DAgger [41] as our imitation learning algorithm and we investigate the performance of the control policy with respect to DAgger iterations. As the DAgger runs increase, the difference in performance between the policies becomes less pronounced as the data aggregation process itself makes all three policies more robust. This provides us with one measure of how effective our embedding method is for data augmentation, particularly in the beginning of imitation.

TABLE I
OUT-OF-DISTRIBUTION EXPERIMENT RESULTS (FLYING VEHICLE)

Training Setup	Interventions	Autonomy(%)	Cross Track Error (m)
Center Camera	3	84.23	2.12
Equiv. Augment. Out-Of-Dist Training	2	89.46	1.39
Equiv. Augment. In-Dist Training	2	89.51	1.26
All Cameras	0	100	0.17

TABLE II
AUGMENTATION LABEL CALCULATIONS (FLYING VEHICLE)

Center Camera	δyaw (rad)	δz (m)
Left Camera	$\delta yaw + 0.03 \text{ rad}$	δz
Right Camera	$\delta yaw - 0.03 \text{ rad}$	δz
Top Camera	δyaw	$\delta z + 0.5m$
Bottom Camera	δyaw	$\delta z - 0.5m$

A. Quadrotor Simulation Experiments

Experimental Setup. For this set of experiments, we use a flying vehicle in the AirSim simulator [42], and also a more complex Drone-Racing environment developed by Microsoft [43], consisting of outdoor three dimensional racing trail marked by rectangular gates. The vehicle is set up with five front facing cameras. For the dataset used in these experiments, each data sample comprises of an image set $I = \langle I_l, I_c, I_r, I_t, I_b \rangle$ corresponding to timestamped images from the left, center, right, top, and bottom cameras. The action tuple $A = \langle \delta yaw, \delta z \rangle$, is the expert control used to navigate the flying vehicle, applied in terms of the relative change in yaw and z at each timestep. The augmented actions for non-center cameras are computed as shown in Table II. Note that clockwise rotation is positive δyaw and downwards direction is positive δz .

Results. Fig. 4 (left) shows the cross track error vs DAgger iterations for three different test scenarios. We see that using the augmentations generated via the equivariant maps for training improves the imitation performance. Fig. 5 overlays the trajectory flown on top of the reference trajectory. The policy trained with equivariant augmentations is able to track better than a policy trained using only the center-camera data. Additionally, as shown in Fig. 6, we observe that the policy trained on our augmentations requires on average 2 fewer interventions when compared to the policy trained without them, for the Drone Racing environment.

Out-Of-Distribution Experiments. We perform another set of experiments where we train the equivariant map on a dataset captured in a separate area of the simulation map, a dense urban environment with buildings and roads, compared to where we test the model, a park with no city structures except a single paved path. As shown in Table I, the policy trained using the generated augmentations has a lower cross track error by 0.73 meters when compared to the policy trained directly on the center-camera data. Furthermore, its cross track error is only 0.13 meters higher than the policy which uses the equivariant map trained on the same area of the simulation map as used for testing.

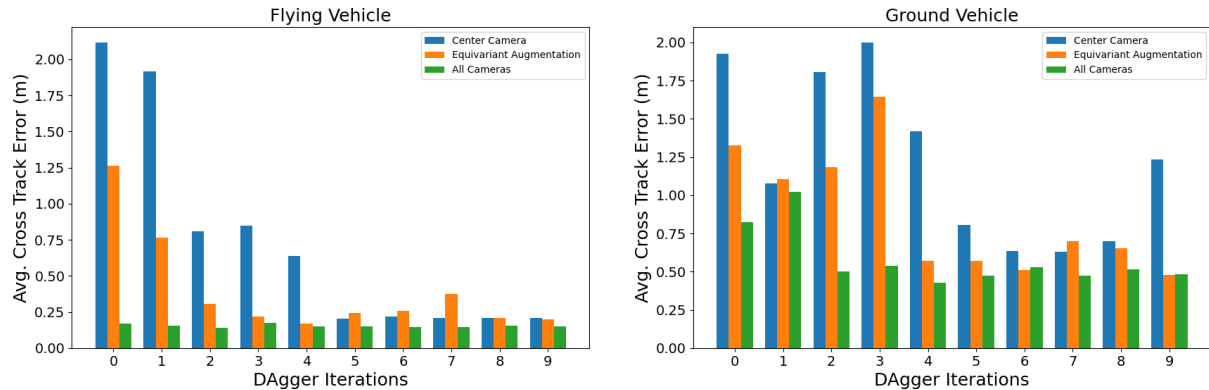


Fig. 4. Cross track error (m) vs DAgger iterations for flying vehicle (left) and ground vehicle (right). It can be seen that network trained with equivariant augmentation lies in between the network trained using center camera and all cameras in terms of performance.

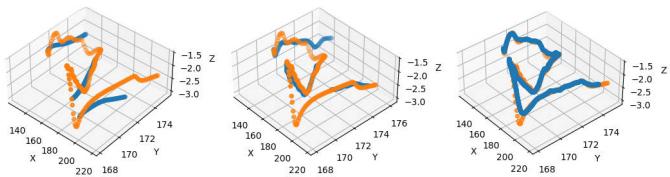


Fig. 5. Reference trajectory (orange) vs flown trajectory (blue). The all-cameras model (right) tracks the trajectory the best, followed by the equivariant-augmentation model (center), followed by the center-camera model (left). Note that missing sections in the flown trajectory are due to an intervention where expert control was used to navigate the vehicle back to the reference trajectory. See IV-A for details.

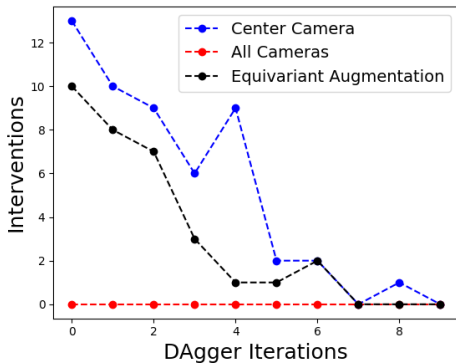


Fig. 6. Number of interventions vs DAgger iterations for the flying vehicle in drone-racing environment. Model trained using equivariant augmentations lies in between the model trained using a single cameras and the model trained using all five cameras.

B. Driving Simulation Experiments

Experimental Setup. For this experiment, we use the Carla simulator [44]. The driving vehicle is set up with three front facing cameras. Each data sample of the dataset consists of an image set $I = \langle I_l, I_c, I_r \rangle$ of timestamped images from the left, center, and right cameras respectively, paired with their corresponding expert control, in this case, the steering angle for the car collected from the Carla driving autopilot. The augmentations for non-center cameras are computed as shown in Table IV.

TABLE III

OUT-OF-DISTRIBUTION EXPERIMENT RESULTS (GROUND VEHICLE)

Training Setup	Interventions	Autonomy (%)	Cross Track Error (m)
Center Camera	8	80.24	4.00
Equiv. Augment. Out-Of-Dist Training	6	87.21	2.66
Equiv. Augment. In-Dist Training	5	88.01	2.61
All Cameras	4	89.97	2.25

TABLE IV

AUGMENTATION LABEL CALCULATIONS (GROUND VEHICLE)

Training Camera	Steering Angle (normalized $[-1, 1]$)
Left Camera	$steering_{cen} + 0.05$
Right Camera	$steering_{cen} - 0.05$

Results. Again we can see in Fig. 4 (right) that the driving policy trained using equivariant augmentations outperforms the policy trained with just the center camera directly, and that this benefit slowly subsides with more DAgger iterations.

Out-Of-Distribution Experiments. As per the flying vehicle experiments, we perform an additional experiment to test generalization of the learned features to novel environments. The experiments so far used a specific path in 'Town10' from the Carla simulator [44] which resembles a densely packed urban environment with lots of city structures. For this experiment however, the equivariance map was trained using a dataset collected on a highway surrounded by green hills and no other city structures. Yet our results as shown in Table III illustrate the significant performance improvement gained by our generated augmentations over the policy trained without them.

Gaussian Noise Augmentation Experiments. We also compare our method to simply adding noise augmentations to the camera images. The noise augmentations consist of zero mean gaussian noise with varying standard deviations $\sigma = \{0.01, 0.05, 0.1, 0.2, 0.3\}$. While noise augmentations do improve navigation performance, our method still performs better in most cases as shown in Fig. 7.

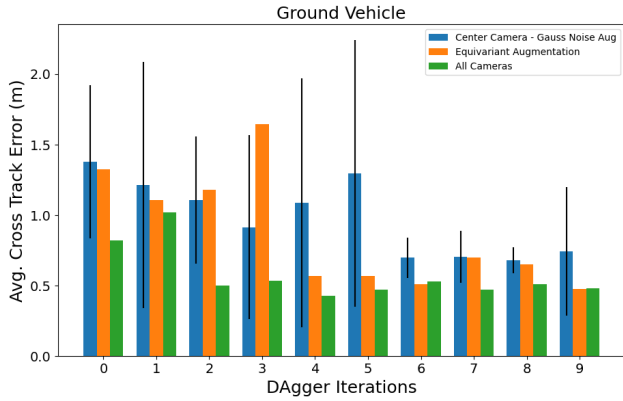


Fig. 7. Average cross track error vs DAGger iterations with center camera network trained with gaussian noise augmentations on the input images.

C. Terrestrial Robot Experiments

Experimental Setup. These experiments were performed using a Clearpath Husky robot equipped with three front facing cameras. The Husky setup is similar to the driving vehicle in Carla simulation as described in Sec. IV-B. Three ZED cameras (used as monocular cameras) are mounted on the front, each 25cm apart. All of the computation was performed on the NVIDIA Jetson Xavier mounted on the bottom. The robot can be seen in Fig. 8.



Fig. 8. Clearpath Husky used for ground robot experiments with three front facing cameras mounted in the front.

Fig. 9 shows the overhead view of the test path used for the experiment. The path is highlighted in red. Three different models are trained to drive the robot along the path—a model trained using data from a single camera, a model trained using data from a single camera with equivariant augmentations, and a model trained using data from all three cameras.

Results. The model trained using only center camera data performs worse than the others, steering off the path 9 times. The model trained on data generated using equivariant maps as well as the model trained using all three cameras both perform significantly better, steering off the path only once. Fig. 9 shows a snapshot of the entire test run for three setups respectively. The locations where an intervention occurred

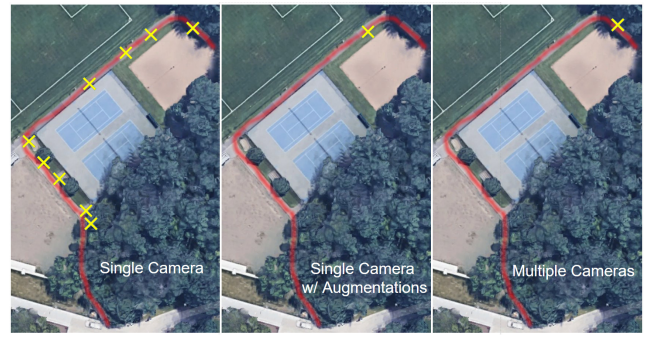


Fig. 9. A bird’s eye view of the testing path for the Husky experiment. Shows the interventions needed while navigating for each policy.

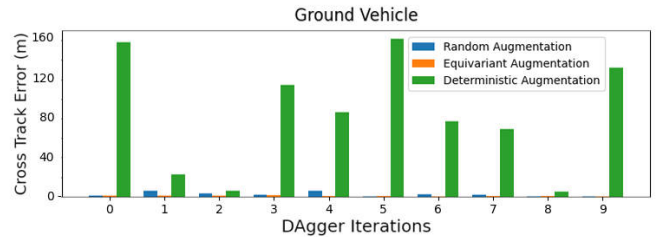


Fig. 10. Cross track error(m) for driving policy trained using a deterministic map and a learned map enforcing equivariant features.

are denoted by a yellow cross.*

D. Learnt vs Non-Learnt Equivariance Maps

In each of the experiments thus far, we modeled the equivariance maps using a neural network that is optimized using expert data. We investigate the necessity of learning these maps by comparing them to a random map that does not enforce any equivariance relation, and a fixed deterministic map (a linear translation in embedding space). Our results as shown in Fig. 10 indicate that both alternative techniques fall short of achieving the performance improvements gained using an equivariance map.

V. CONCLUSION

In this paper, we investigate the use of equivariant maps for visual navigation by mobile robots. By training a neural network model that learns image features which are equivariant, we can predict the latent representations of images from viewpoints nearby to the one observed. Our results indicate that by augmenting the training dataset with these representations, one can significantly improve navigation performance in a variety of settings as demonstrated by our simulation experiments involving flying and ground vehicles. Additionally, through our ground robot experiments over a 500m path, we showed that the benefits of our method also transfer over to real world settings.

*The video from this experiment can be found here: <https://youtu.be/5g4Kg3-YWvA>

REFERENCES

- [1] J. Cohen, E. Rosenfeld, and Z. Kolter, "Certified adversarial robustness via randomized smoothing," in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 09–15 Jun 2019, pp. 1310–1320.
- [2] D. Jayaraman and K. Grauman, "Learning image representations tied to ego-motion," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [3] K. Lenc and A. Vedaldi, "Understanding image representations by measuring their equivariance and equivalence," 2015.
- [4] U. Schmidt and S. Roth, "Learning rotation-aware features: From invariant priors to equivariant descriptors," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 2050–2057.
- [5] Y. Zhao, T. Birdal, J. E. Lenssen, E. Menegatti, L. Guibas, and F. Tombari, "Quaternion equivariant capsule networks for 3d point clouds," *ArXiv*, vol. abs/1912.12098, 2019.
- [6] T. Tuytelaars and K. Mikolajczyk, "Local invariant feature detectors: A survey," *Found. Trends. Comput. Graph. Vis.*, vol. 3, no. 3, p. 177–280, Jul. 2008. [Online]. Available: <https://doi.org/10.1561/06000000017>
- [7] J. J. Kivinen and C. K. I. Williams, "Transformation equivariant boltzmann machines," in *Artificial Neural Networks and Machine Learning – ICANN 2011*, T. Honkela, W. Duch, M. Girolami, and S. Kaski, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 1–9.
- [8] M. Bilovs and S. Gunnemann, "Equivariant normalizing flows for point processes and sets," 2020.
- [9] D. J. Rezende, S. Racanière, I. Higgins, and P. Tóth, "Equivariant hamiltonian flows," *ArXiv*, vol. abs/1909.13739, 2019.
- [10] M. Weiler, M. Geiger, M. Welling, W. Boomsma, and T. Cohen, "3d steerable cnns: Learning rotationally equivariant features in volumetric data," in *NeurIPS*, 2018.
- [11] R. Spezialetti, S. Salti, and L. Stefano, "Learning an effective equivariant 3d descriptor without supervision," *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 6400–6409, 2019.
- [12] F. Fuchs, D. E. Worrall, V. Fischer, and M. Welling, "Se(3)-transformers: 3d roto-translation equivariant attention networks," *ArXiv*, vol. abs/2006.10503, 2020.
- [13] R. Wang, M. Albooyeh, and S. Ravanbakhsh, "Equivariant maps for hierarchical structures," *ArXiv*, vol. abs/2006.03627, 2020.
- [14] D. Worrall and G. Brostow, "Cubenet: Equivariance to 3d rotation and translation," in *Computer Vision – ECCV 2018*, V. Ferrari, M. Hebert, C. Sminchisescu, and Y. Weiss, Eds. Cham: Springer International Publishing, 2018, pp. 585–602.
- [15] D. E. Worrall, S. J. Garbin, D. Turmukhambetov, and G. J. Brostow, "Harmonic networks: Deep translation and rotation equivariance," 2017.
- [16] T. S. Cohen, M. Weiler, B. Kicanaoglu, and M. Welling, "Gauge equivariant convolutional networks and the icosahedral cnn," 2019.
- [17] K. S. Tai, P. Bailis, and G. Valiant, "Equivariant transformer networks," *ArXiv*, vol. abs/1901.11399, 2019.
- [18] T. Cohen, M. Geiger, and M. Weiler, "A general theory of equivariant cnns on homogeneous spaces," *ArXiv*, vol. abs/1811.02017, 2019.
- [19] M. Weiler, F. A. Hamprecht, and M. Storath, "Learning steerable filters for rotation equivariant cnns," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 849–858.
- [20] C. Esteves, Y. Xu, C. Allec-Blanchette, and K. Daniilidis, "Equivariant multi-view networks," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 1568–1577.
- [21] A. Kanazaki, Y. Matsushita, and Y. Nishida, "Rotationnet: Joint object categorization and pose estimation using multiviews from unsupervised viewpoints," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5010–5019.
- [22] Y. Feng, Z. Zhang, X. Zhao, R. Ji, and Y. Gao, "Gvcnn: Group-view convolutional neural networks for 3d shape recognition," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 264–272.
- [23] C. R. Qi, H. Su, M. Nießner, A. Dai, M. Yan, and L. Guibas, "Volumetric and multi-view cnns for object classification on 3d data," in *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2016.
- [24] H. You, Y. Feng, R. Ji, and Y. Gao, "Pvnet: A joint convolutional network of point cloud and multi-view for 3d shape recognition," 2018.
- [25] E. van der Pol, T. Kipf, F. A. Oliehoek, and M. Welling, "Plannable approximations to mdp homomorphisms: Equivariance under actions," in *AAMAS*, 2020.
- [26] A. J. Ratner, H. R. Ehrenberg, Z. Hussain, J. Dunnmon, and C. Ré, "Learning to compose domain-specific transformations for data augmentation," 2017.
- [27] X. Zhang, Q. Wang, J. Zhang, and Z. Zhong, "Adversarial autoaugmentation," 2019.
- [28] T. Dao, A. Gu, A. J. Ratner, V. Smith, C. D. Sa, and C. Ré, "A kernel theory of modern data augmentation," 2019.
- [29] R. Raileanu, M. Goldstein, D. Yarats, I. Kostrikov, and R. Fergus, "Automatic data augmentation for generalization in deep reinforcement learning," 2020.
- [30] R. Garg, V. K. B.G., G. Carneiro, and I. Reid, "Unsupervised cnn for single view depth estimation: Geometry to the rescue," in *Computer Vision – ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds. Cham: Springer International Publishing, 2016, pp. 740–756.
- [31] A. Loquercio, A. I. Maqueda, C. R. D. Blanco, and D. Scaramuzza, "Dronet: Learning to fly by driving," *IEEE Robotics and Automation Letters*, 2018.
- [32] U. Müller, J. Ben, E. Cosatto, B. Flepp, and Y. L. Cun, "Off-road obstacle avoidance through end-to-end learning," in *Advances in Neural Information Processing Systems 18*. MIT Press, 2006, pp. 739–746.
- [33] D. Pomerleau, "Neural network perception for mobile robot guidance," Ph.D. dissertation, Carnegie Mellon University, 1993.
- [34] D. A. Pomerleau, "Alvinn: An autonomous land vehicle in a neural network," in *Advances in Neural Information Processing Systems 1*, D. S. Touretzky, Ed., 1989, pp. 305–313.
- [35] S. Ross, G. J. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *International Conference on Artificial Intelligence and Statistics, AISTATS*, 2011, pp. 627–635.
- [36] A. Giusti, J. Guzzi, D. C. Cireşan, F. L. He, J. P. Rodriguez, F. Fontana, M. Faessler, C. Forster, J. Schmidhuber, G. D. Caro, D. Scaramuzza, and L. M. Gambardella, "A machine learning approach to visual perception of forest trails for mobile robots," *IEEE Robotics and Automation Letters*, vol. 1, no. 2, pp. 661–667, July 2016.
- [37] M. Bojarski, D. D. Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Müller, J. Zhang, X. Zhang, J. Zhao, and K. Zieba, "End to end learning for self-driving cars," *CoRR*, vol. abs/1604.07316, 2016. [Online]. Available: <http://arxiv.org/abs/1604.07316>
- [38] F. Codevilla, M. Müller, A. Dosovitskiy, A. López, and V. Koltun, "End-to-end driving via conditional imitation learning," *CoRR*, vol. abs/1710.02410, 2017.
- [39] M. Laskey, J. Lee, R. Fox, A. Dragan, and K. Goldberg, "Dart: Noise injection for robust imitation learning," in *Proceedings of the 1st Annual Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, vol. 78. PMLR, 13–15 Nov 2017, pp. 143–156.
- [40] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015. Conference Track Proceedings*, Y. Bengio and Y. LeCun, Eds., 2015. [Online]. Available: <http://arxiv.org/abs/1412.6980>
- [41] S. Ross, G. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, 2011, pp. 627–635.
- [42] S. Shah, D. Dey, C. Lovett, and A. Kapoor, "Airsim: High-fidelity visual and physical simulation for autonomous vehicles," in *Field and Service Robotics*, 2017. [Online]. Available: <https://arxiv.org/abs/1705.05065>
- [43] R. Madaan, N. Gyde, S. Vemprala, M. Brown, K. Nagami, T. Taubner, E. Cristofalo, D. Scaramuzza, M. Schwager, and A. Kapoor, "Airsim drone racing lab," *arXiv preprint arXiv:2003.05654*, 2020.
- [44] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in *Proceedings of the 1st Annual Conference on Robot Learning*, 2017, pp. 1–16.