# Learning autonomous driving from aerial imagery

Varun Murali[1], Guy, Rosman[2], Sertac, Karaman[1] and Daniela Rus[3]

*Abstract*— In this work, we consider the problem of learning end to end perception to control for ground vehicles solely from aerial imagery. Photogrammetric simulators allow the synthesis of novel views through the transformation of pre-generated assets into novel views. However, they have a large setup cost, require careful collection of data and often human effort to create usable simulators. We use a Neural Radiance Field (NeRF) as an intermediate representation to synthesize novel views from the point of view of a ground vehicle. These novel viewpoints can then be used for several downstream autonomous navigation applications. In this work, we demonstrate the utility of novel view synthesis though the application of training a policy for end to end learning from images and depth data. In a traditional real to sim to real framework, the collected data would be transformed into a visual simulator which could then be used to generate novel views. In contrast, using a NeRF allows a compact representation and the ability to optimize over the parameters of the visual simulator as more data is gathered in the environment. We demonstrate the efficacy of our method in a custom built mini-city environment through the deployment of imitation policies on robotic cars. We additionally consider the task of place localization and demonstrate that our method is able to relocalize the car in the real world.

## I. INTRODUCTION

With the rising popularity of self-driving cars and the need for high definition maps for navigation, it is becoming ever necessary to gather data for training autonomous driving on urban roads. This requires collecting continuous driving data and is limited to regions which are heavily trafficked or allowable by the rules of the road. On the other hand, aerial imagery captures larger regions from a different viewpoint which can be transformed into the viewpoint of a ground robot. This allows the capture of edge cases or visual information that might be otherwise occluded cumbersome to represent in visual reconstructions or simulators. Visual simulators have also grown increasingly in capability with the advancement of real time ray tracing solutions such as Nvidia DLSS and the availability of commercial photogrammetric tools such as RealityCapture [1], Matterport [2].

Aerial images in contrast are readily available in the form of satellite imagery and are less cumbersome to collect since they are able to map multiple streets simultaneously. In this work, we operate under the real to sim to real framework where we transform images into novel viewpoints for data generation through the intermediate representation. Since this representation is also in itself learnt, we are able to store the representation in a compact way and optimize over parameters such as lighting, and locations of objects where new
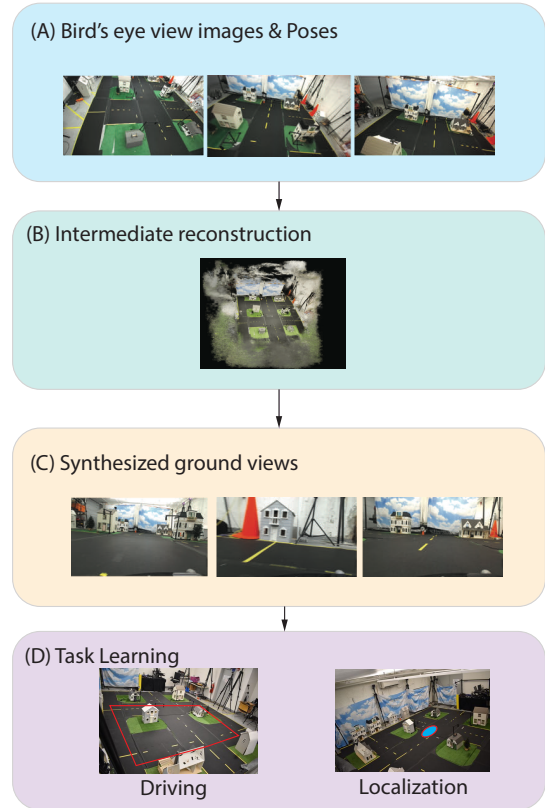
Fig. 1: An overview of the method presented in this work. First, we assume that we are given aerial images and their corresponding poses (A). We then reconstruct a photogrammetric model (B) of the world using BEV images and use this intermediate representation which can be used to query ground robot views for desired poses (C). We use the synthesize images to learn policies that can be directly deployed on real vehicles (D).

data is available. In contrast, a traditional simulator (built through photogrammetry such as FlightGoggles [3]) would require the transfer of imagery into visual assets that can then be transformed into novel views but updating it would require the generation of new assets or creation of a new database for photogrammetry. Photogrammetric simulators are then able to output various visual modalities including the camera imagery, depth imagery and other representations such as semantic segmentation and surface normals. These can then be used to augment the training data available to train autonomous agents (such as in VISTA [4]) or allow

agents to learn policies through trial in simulation (such as in AI habitat [5]). Such simulations allow us to consider data that would not otherwise be available in traditional human collected data, *e.g.*driving off-road or near accident scenarios. Generative simulation [6] is already playing a role in robotics applications across multiple domains. Simulation engines are increasingly used to combine photogrammetric assets and generated ones to train machine learning algorithms but this requires a lot of data and engineering effort. It is also unclear how to capture the day to day changes in the environment with high fidelity in such a simulation.

Learning representations of the physical world also alleviates the need to reason about occlusions during the synthesis of (depth and color) images from the reconstruction. In a traditional setting, the simulator would hold a object mesh of the physical environment and then project rays and reason about occlusions from a certain viewpoint to render the scene onto camera space. In contrast, a neural radiance field learns both the color along the ray and the density along the ray implicitly learning the occupancy map. For this work, we consider learning from these representations two problems in autonomous driving: (i) visual re-localization and (ii) end-to-end driving. Visual (re-)localization is a challenging problem in infrastructure denied navigation. Most driving applications require localization to downstream tasks. End-to-end driving on the other hand attempts to implictly learn the localization. While, these tasks are related in that localization could aid a motion planning algorithm to follow trajectories by considering obstacles in the environment, in this work the second task is learnt completely independently to demonstrate the possibility of implicitly reasoning about plans from only onboard imagery. We stress here that the applications are cross-domain tasks i.e. taking BEV images as input and estimating pose or desired actions in a ground robot's view. We also stress that while we use NeRF as the intermediate perceptual representation for cross-view tasking, the methods presented in this work are general enough to be extended to other methods such as Gaussian Splatting [7]. Fig 1 outlines the proposed components of this work.

In summary, in this work we represent the transformation of a collection of images from a bird's eye view (BEV) into a ground view (rgb imagery) using a neural radiance field as an intermediate representation. We then use the generated ground views to train an imitation policy using a pure-pursuit controller to learn a mapping from the robot's onboard camera to a steering and velocity command. We demonstrate the utility of our method through real world experiments on one-tenth scale cars in a miniature city environment by transferring the learned policies directly onto the vehicles.

## II. RELATED WORK

Our work intersects with several topics of research in robotics and robotic perception. In particular, we leverage ideas from multi-view stereo from aerial images and novel view synthesis from neural rendering fields. We also leverage ideas from reasoning across varied viewpoints and photogrammetric real-to-sim-to-real pipelines which generate photorealistic imagery to train machine learning algorithms. Finally, we demonstrate the efficacy of our approach through the tasks of (i) visuo-motor policy learning and (ii) visual place recoginition.

### A. Multi-View reconstruction from aerial images

Several works have considered the problem of generating dense reconstruction from satellite and aerial images such as [8], [9], [10], [11], [12], [13], [14]. These works focus on matching features across the images and aligning them to minimize a photometric error to recover a dense surface model. Another set of related works involved ground-and-aerial navigation and planning [15], [16]. Since these works usually involve cumbersome feature detection, matching and global bundle adjustment to generate the model they are computationally expensive to compute and would require an additional meshing procedure to store the generated surface models. Additional reasoning is often needed to cull occluded points and vary parameters like lighting to produce realistic images that are actionable and artifact-free. Leveraging neural radiance fields as we do in this work help alleviate some of these challenges.

### B. Novel View Synthesis

Neural radiance fields were first presented by Mildenhall et al. [17] to learn the mapping of a projected ray to the colors and the log density along the ray. Since the initial appearance of the neural radiance fields, several works has considered improvements such as relighting the scene [18], relaxing the requirement of accurate poses [19], extensions to style transfer and near-real time field generation [20]. [21], [20] Recent work also considers the idea of regularizing the predicted densities in order to improve the quality of the depth maps [22]. Multi-view NeRFs have also been proposed to better utilize nearby reference views at the time of inference [23]. City scale reconstructions using NeRFs have also been considered [24]. This effectively stores the problem of storing large scenes but is trained on similar views as the rendered views. 3D Gaussian splatting based approaches have also been used as an alternative to neural radiance fields to learn dense representations from multi-view camera imagery [7]. In this work, we use NeRF as a backbone representation to store the perceptual information but stress that the proposed method is easily extensible to other representations.

### C. Representations for Cross-view Embodied Reasoning

The problem of learning to act in the world based on a different view of the scene is been the focus of several research threads. In trajectory generation and forecasting, several works have looked at birds-eye-views as an input for acting. [25] have learned a latent map representation from BEV images and embedded local cues for acting into it. Adamkiewicz *et al.*[26] have leveraged NeRFs to create an occupancy estimate and planned in the resulting volume. For localizing agents, [19] have demonstrated how to correct camera poses along with reconstructing the NeRF,

whereas Moreau *et al.*[27] proposed how to leverage NeRFs to enhance localization. Semantic representations for cross-view localization have also been studied [28], [29] where the geometric composition of semantically important objects in the scene are used to register aerial views to ground views.

### D. Photogrammetric simulators

Simulators such as FlightGoggles [3] aim to solve the sim to real to sim by first generating the photogrammetric assets through software and human effort and then feeding it to a game engine such as Unity3D. In this scenario, the assets are "baked" into the simulator and the approximate ray tracing methods relevant in the computer graphics community can be used to relight the scene and synthesize images. However, while such simulators are capable of fast generation of imagery, the assets themselves require high level of detail to store and are often expensive in terms of memory. Any changes to the physical environment would also require repeating the asset capture process (depending on the granularity of the change). Simulators such as VISTA alleviate the photogrammetric asset creation by generating only local transformations from pre-existing datasets. While this allows for a diversity in generating training data, it does not allow the introduction of imagery from a drastically different view point. Photogrammetric simulation also often takes a lot of time to generate and requires capturing a lot of imagery for good performance.

### E. Visuo-motor policies

Prior work also considers learning policies that ingest vision and directly regress the desired robot control. [30] address the question of training perception and control jointly and demonstrate the ability to learn manipulation tasks directly from imagery. [31] use imitation learning to learn end-to-end driving given on-board imagery from the car and noisy localization estimate. [32], [33], [34] demonstrate learning end-to-end driving from a large amount of driving video data. [35] use NeRFs to generate simulation environments from scenes by collecting data from a phone and then learn locmotion policies in this simulation. Although the general application and idea is similar, in this work we utilize images from different domains i.e. BEV to train the NeRF and ground views for the simulation.

### F. Visual place recognition

Work has also considered the problem of localization from single images given a visual library or map [36]. For instance, [37] learn to regress the camera position and orientation given a sequence of images and poses from the same viewpoint and same region. Newer work has also considered learning or learnt feature representations to hash maps and estimate the pose from images [38]. In contrast to prior work, we directly leverage aerial imagery to perform cross view tasks that are relevant to autonomous driving. Specifically, we demonstrate the ability to localize the ego-vehicle and perform visuo-motor control onboard ground vehicles with only imagery from an aerial vehicle.

## III. METHOD

In this section, we first detail the problem setup. We then explain the procedure for generating the intermediate NeRF based simulation. Then, we present two two learning problems that are relevant for autonomous driving.

### A. Problem Setup

Let's denote images $I$ with the superscript $A$ representing a bird's eye view and the superscript $G$ represent a ground view. The transform $_w^r T$ represents the transform from the robot to the world reference frame and $\mathcal{K}$ represent the camera matrix. Given a sequence of images rendered from the point of view of the aerial robot represented by $I_{1,\cdots,n}^A$ and their locations $_w^r T_{1,\ldots,n}$ we wish to find a policy $\phi$ that maps an image taken from the point of view of the ground vehicle $I^G$ to a control command comprising the steering $\omega$ and velocity $v$ command. This requires the implicit construction of the scene represented by $c, \sigma$ where $c$ represents a function mapping the color to a physical location in the world and $\sigma$ represents a function mapping the probability of the location being occupied. Figure 2 shows the full pipeline of our proposed approach.

### B. Reconstruction

**Model.** Given the sequence of images $I_{1,\cdots,n}^A$, we learn a neural radiance field using a multi-resolution hash approach [20]. As a backbone network for the NeRF, we use the method presented by [20] to learn two concatenated multilayer perceptron. The first learns the density i.e. a mapping $\sigma = \Psi(x; \theta)$ from the encoded position and direction of a projected ray to the density at that location. The second learns the color and allowing for view-dependent color variation. We use an exponential moving average optimizer with a learning rate of $1e^{-3}$ and $\beta_1 = 0.9$ and $\beta_2 = 0.99$ for 5000 steps.

**Adding road priors.** Since our method is reliant on cross view transfers and several occlusions exist in the aerial views, we enforce an additional prior on the depth and color $c_r$ of the road networks. We assume here that the location of the road network is available to us apriori. We randomly include rays that are normal to the road surface with the predefined color to the training batch. For each ray that is normal to the road surface, we first sample a position along that ray and include the position, view direction and pre-defined color into the training buffer. We assume that the distance between the sampled ray position and the robot is equal to the height above ground of the aerial vehicle.

**Losses.** For each image in the training set we compute the loss of between the reconstructed image by projecting rays from the image location $T_w^I$ and the training images using a photometric error metric. The road prior rays also contribute to the loss by computing the error between the predicted color and the desired road color. Additionally, since the poses of the cameras are noisy we optimize over the poses of the cameras by computing the photometric error of projecting a nearby view onto the desired view.
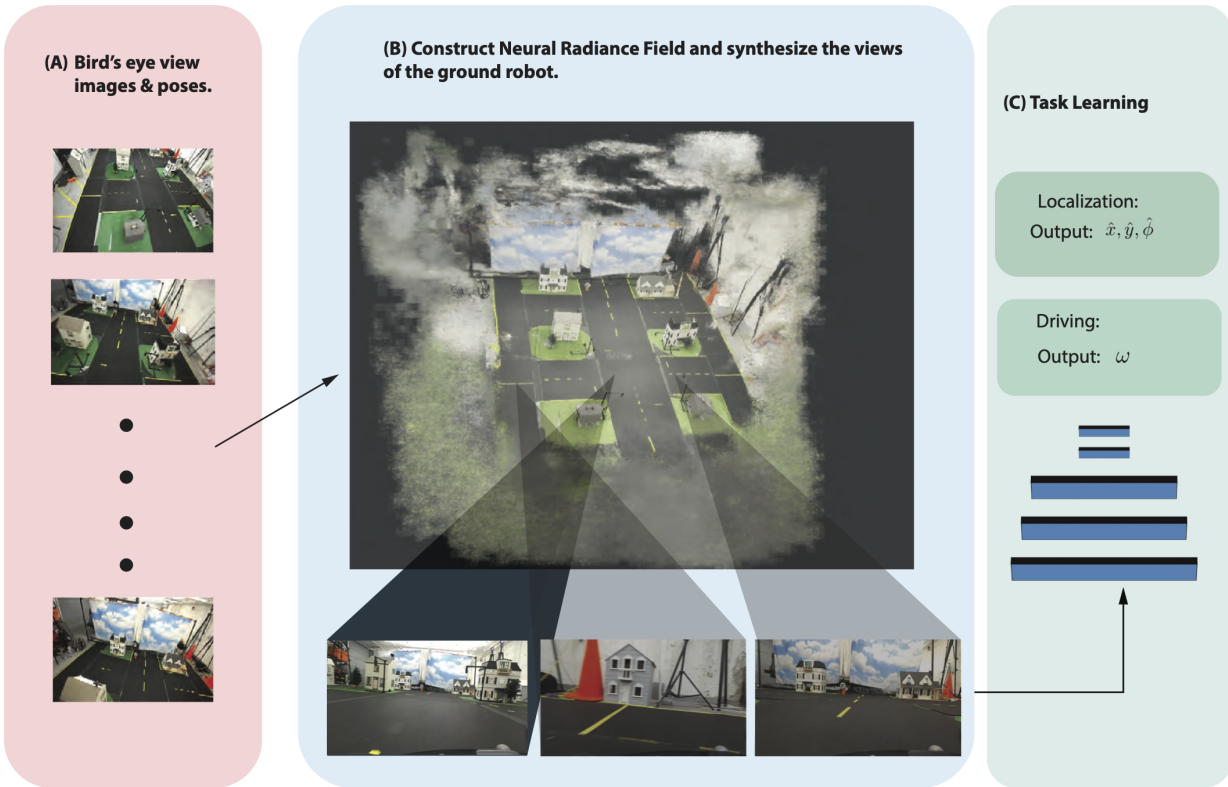
Fig. 2: An overview of our proposed method. (A) We assume that are we are given images and their corresponding poses from the bird's eye view. (B) We leverage a neural render field to compactly represent the density and color of the desired scene. The NeRF can be queried for poses along the road network. (C) We then learn two task relevant to autonomous driving: (i) visual localization and (ii) end-to-end driving.

**Algorithm.** In summary, the algorithm used for reconstruction from the bird's eye view image sequence is is presented in Algorithm 1. The input to the algorithm is the sequence of images and this stage, we produce a trained neural network that can output density and color given a ray. The function ray takes a pixel and pose and returns the ray between them.

*C. Data generation and models*

**Image Generation.** Given the reconstruction, we can synthesize images at novel locations given a query pose. We project rays consistent with the camera matrix $\mathcal{K}$ from each query location to get a set of RGB values and densities along the ray. Assuming the density is a log probability of the ray producing a "hit", we can then sum the product of the predicted colors and the densities to get a color value for that pixel. To generate adequate training data, we sample poses along the known road network. We take care to sample poses that would also violate driving rules (*e.g.* driving on the wrong side of the road) to provide data augmentation. This is an advantage of using an intermediate reconstruction or "simulation" since we can generate views that would not normally be present in collected driving data. Prior work has already demonstrated the importance of this [4].

**Control generation.** In this work, we use the idea of pure-pursuit control as formulated in [39]. However, our method

---

**Algorithm 1** Procedure for reconstruction.

**Require:** ${}_w^r T_{1,\cdots,n}, I_{1,\cdots,n}^A, N$
  **for** $x = 1 \rightarrow N$ **do**
    **for** $k = 1 \rightarrow n$ **do**
      **for** pixel in $I_k^G$ **do**
        $x, \theta, c = \text{ray}({}_w^r T_k, \text{pixel})$
        Buffer $\leftarrow (x, \theta, c)$
      **end for**
      **for** pixel in road network **do**
        $T \leftarrow$ random sample $\perp$ to road
        $x, \theta, c_r = \text{ray}(T, \text{pixel})$
        Buffer $\leftarrow (x, \theta, c_r)$
      **end for**
      Shuffle(Buffer)
      TrainNerf(Buffer)
    **end for**
  **end for**

---

is generic and can admit other types of control schemes in the policy training if desired. In the pure-pursuit algorithm a target path is known, a point on this path some fixed distance ahead of the vehicle is targeted, and a heading command is derived from the car's position and dimensions such that it will intersect this point within an acceptable margin. For

this control scheme, we assume that the output velocity is constant. An external rule based controller is used to change the velocities around obstacles. We also assume that the controller acts as a low pass filter on the steering commands to prevent over and under steering due to noisy inputs.

---

**Algorithm 2** Procedure for task learning.

---

**Require:** $^r_w T_{1,\cdots,n}, N, task$
  **for** $x = 1 \rightarrow N$ **do**
    **for** $k = 1 \rightarrow n - 1$ **do**
      $^r_w T_k \leftarrow ^r_w T_k + perturb$
      **for** pixel in $I^G_k$ **do**
        $x, \phi = \text{projectRay}(^r_w T_k)$
        pixel $\leftarrow c(x, \phi)$
      **end for**
      **if** task = visuo-motor **then**
        $\omega, v \leftarrow \text{computeControl}(^r_w T_k, ^r_w T_{k+1})$
        Buffer $\leftarrow I^G_k, (\omega, v))$
      **else**
        Buffer $\leftarrow I^G_k, (^w_r T))$
      **end if**
    **end for**
    Shuffle(Buffer)
    TrainModel(Buffer)
  **end for**

---

**Model.** For both tasks under consideration, we use a five layer convolutional with 8 filters neural network with batch normalization on each layer and a rectified linear unit for activation to encode the images. The output of the convolutional neural network is fed through two fully connected layer with 64 hidden units each to output the steering command for the driving task and to output the position and yaw for the localization task. We implement the neural network in the pytorch framework [40] with the ADAM optimizer [41] with a learning rate of $1e^{-2}$ for 100 epochs. We stress that we intend here to keep the model size small to be feasibly deployed on the embedded platform on our robots but emperically we observe that these networks elicit the desired performance.

### D. Tasks

To demonstrate the cross-view transfer and it's applicability to autonomous driving, we choose the tasks of visual pose estimation and visuo-motor policy learning but stress that the method could be applicable to other tasks. We use the models described in the previous subsection. The models are trained independently per-task and per-environment and do not share any parameters. We also highlight that the second task is end-to-end and implicitly reasons about localization.

**Visual pose estimation.** The first task we consider is that of estimating the robot pose given a single image $I^G$ and the relative transformation between the camera and the robot $^c_r T$. For this task, we assume that the configuration of the environment is static and does not change between when the bird's eye view imagery was collected and the policy was deployed onboard the ground vehicle. In this case, we train

the neural network on all feasible locations of the ground robot for 100 epochs using Algorithm 2. To train this model, we compute the loss using an L2 loss between the ground truth image pose and the predicted pose. Concretely, the goal here is to learn directly the desired steering angle of the car given the current camera image.

**Visuo-Motor Policy Learning.** The next task under consideration is that of learning steering commands given the onboard imagery of the robot. We assume that we are given a predefined trajectory that we would like to follow in the world frame. Given the trajectory in the world frame, we sample a sequence of positions and orientations along that path and generate images from the reconstruction from the point of view of the car at those locations. The procedure to train the model for 100 epochs is given in Algorithm 2. To train this model, we compute the loss using an L1 loss between the desired steering and the predicted steering command. The goal of this model is to learn directly the desired steering angle of the car given the current camera image.

## IV. EXPERIMENTS

In this section, we detail the experimental setup and experiments to validate our method. First, we discuss the experimental platform used for validation, the data collection pipeline and then the experiments for the two tasks: (i) localization and (ii) end-to-end driving. We assume in this section, that a neural network is trained for each task and each configuration of the environment.

### A. Mini-city experimental platform

We use the mini-city experimental platform described in [42] and visualized in Figure 3. We deploy the learned policies onboard a custom built racecar built on the MIT racecar platform [43]. The racecar carries a Jetson Xavier NX which is used to run the policies and a ZED 2i camera to acquire images and inertial measurements. We use a Parrot Bebop 2 drone to collect images in the mini-city. We employ 12 Optitrack Prime 41W cameras to estimate the poses of robots inside the minicity for associating the poses from the bird's eye view and for ground truth for the driving tasks. For each experiment, we collect approximately 100 images from the aerial view to generate the NeRF.

### B. Qualitative comparisons

First, we evaluate the visual quality of the cross view transfer using the NeRF as an intermediate representation. We desire to study the effect of using the road priors on the intermediate reconstruction. As can be seen in Figure 5, we can see that the rendered images from the point of view of the ground robot captures details close to the ground and reasonable captures the geometry of environmental features such as houses.

### C. Visual Localization

**Setup.** In this experiment, we measure the re-localization error using the localization policy. The ground robot is
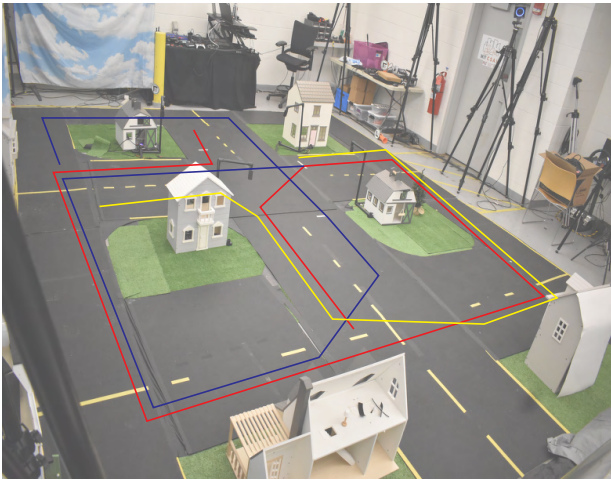
Fig. 3: The figure shows the desired trajectories in the minicity environment. The first trajectory is shown in red, the second in blue, and the third in yellow.



Fig. 4: Experimental setup used. The figure shows an example configuration of the houses, the road network, the motion capture system and the Parrot Bebop drone used to collect data.



Fig. 5: Rendered image from the NeRF from the point of view of the ground robot with ground priors (top) and without(bottom).

baseline and our method is similar in the number of poses sampled along the trajectory.

**Results.** The results of the experiment are summarized in Table I. As can be seen in the table, our method outperforms the baseline when the inertial measurements are added to the localization even though no ground view imagery is used while training. Our method is also comparable to the baseline when then inertial measurements are not included.

manually driven around the environment while running the localization network onboard. The imagery from the ZED camera on the ground robot is fed into the network and the output pose is recorded. We also recorded a filtered estimate that fuses the output pose from the neural network and the inertial measurements on the robot. The different configurations of the environment used for the visual localization task is shown in Figure 6.

**Metrics.** For this experiment, we show the root mean square error of estimating the poses along the trajectory for the position ($\|x - x_{ref}\|_2$) and orientation ($\|\psi - \psi_{ref}\|$) with respect to the ground truth collected using motion capture.

**Baseline.** As a baseline, we train our network on data collected from manually driving the ground robot around the environment. The number of training images for both the

| Method | $\|x - x_{ref}\|_2$ [m] | $\|\psi - \psi_{ref}\|$ [deg] |
|---|---|---|
| **Scenario 1** | | |
| Baseline | 0.13 | 5.31 |
| Ours | 0.17 | 7.82 |
| Ours + IMU | 0.11 | 2.26 |
| **Scenario 2** | | |
| Baseline | 0.15 | 7.89 |
| Ours | 0.21 | 8.61 |
| Ours + IMU | 0.12 | 1.59 |
| **Scenario 3** | | |
| Baseline | 0.12 | 3.45 |
| Ours | 0.18 | 7.44 |
| Ours + IMU | 0.09 | 1.12 |

TABLE I: The root mean square estimation error of the position and yaw angle of the ground robot with respect to the motion capture system.

Fig. 6: The different configurations of the environment used for the visual localization task.

## D. End to End Control

**Setup.** In this experiment, we assume that are we given a specific set of desired poses for the ground robot and we train the end-to-end driving network for that sequence of poses. The desired trajectories are shown in Figure 3. We allow the network to specify the steering angle given the onboard imagery while keeping the nominal speed constant. For this experiment, we record the pose of the driven trajectory and the ground truth from the motion capture system.

**Metrics.** We measure the number of interventions required while deploying the learnt policies and the root mean square error for pose ($\|x - |x_{ref}\|_2$) and orientation ($\|\psi - \psi_{ref}\|$) from the pre-defined trajectory.

**Imitation.** The first baseline that we implement is pure imitation learning from collected datasets from the point of view of the ground robot. We collect teleoperated data in the minicity and train the imitation policy using the collected dataset and test the performance by deploying the policy on the robot. Since there is no domain shift, we expect this to perform reasonably with no need for additional sensors.

**Dataset + View Synthesis from ground view.** In the second baseline, we generate a NeRF from the collected datasets and add additional data by perturbing the poses of the robot in the dataset. We include the predicted images from the NeRF for the perturbed path and the required control to bring the robot back to the nominal path in the training set [32]. For all methods, we use the same neural network with the only difference being the input data for fair comparison. This would yield a similar performance and setup as [35] since the domain of the training data would be similar to the task data.

**Results.** The results of this experiment are summarized in Table II. As can be seen in the table, our method outperforms the baselines when no augmentation is used and outperforms the ground view augmentation when also augmenting with the bird's eye view. The results also show that our method requires less interventions than the baseline.

## V. Conclusion

In summary, we have presented an innovative approach for teaching ground robots visuo-motor policies and localization, using solely bird's eye view images and their corresponding

| Method | $\|x - x_{ref}\|_2$ [m] | $\|\psi - \psi_{ref}\|$ [deg] | Interventions |
|---|---|---|---|
| **Trajectory 1** | | | |
| Imit | 0.39 | 6.71 | 10 |
| Imit + GV | 0.23 | 3.89 | 3 |
| Ours | 0.36 | 7.39 | 5 |
| Ours + GV | **0.11** | **2.21** | **1** |
| **Trajectory 2** | | | |
| Imit | 0.27 | 13.8 | 5 |
| Imit + GV | 0.14 | 7.29 | 1 |
| Ours | 0.20 | 7.47 | 2 |
| Ours + GV | **0.09** | **3.73** | **0** |
| **Trajectory 3** | | | |
| Imit | 0.34 | 14.16 | 8 |
| Imit + GV | 0.24 | 10.60 | 2 |
| Ours | 0.16 | 10.75 | 3 |
| Ours + GV | **0.12** | **3.57** | **1** |

TABLE II: The root mean square estimation error of the pose and the number of interventions during execution. Imit: Imitation from the user driven trajectories only. GV: Adding ground view augmentation through the NeRF.

poses as input. This methodology was successfully applied and tested on a 1/10th scale environment and car within the minicity platform. Our findings suggest the potential to enable the acquisition of autonomous driving capabilities with a limited quantity of aerial imagery, in contrast to the conventional approach that relies on extensive datasets. Our results also show that our method can also closely approximate the results of using imagery from the same domain alleviating the necessity of data from the same view. In future work we aim to extend this methodology to full-scale vehicles, operating in large outdoor settings, by harnessing the power of open source satellite data. We can also consider introducing other functions such as dynamic actors into the view synthesis using a time-dependent nerf.

## References

[1] "RealityCapture - 3d Models from Photos and/or Laser Scans," @mischttps://www.capturingreality.com, howpublished = https://www.capturingreality.com, note = Accessed: 2023-09-15 , accessed: 2023-09-15.

[2] S. K. Ramakrishnan, A. Gokaslan, E. Wijmans, O. Maksymets, A. Clegg, J. Turner, E. Undersander, W. Galuba, A. Westbury, A. X. Chang, M. Savva, Y. Zhao, and D. Batra, "Habitat-Matterport 3D dataset (HM3D): 1000 large-scale 3D environments for embodied AI," Sept. 2021.

[3] W. Guerra, E. Tal, V. Murali, G. Ryou, and S. Karaman, "Flightgoggles: Photorealistic sensor simulation for perception-driven robotics

using photogrammetry and virtual reality," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 6941–6948.

[4] A. Amini, I. Gilitschenski, J. Phillips, J. Moseyko, R. Banerjee, S. Karaman, and D. Rus, "Learning robust control policies for end-to-end autonomous driving from data-driven simulation," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1143–1150, 2020.

[5] M. Savva, A. Kadian, O. Maksymets, Y. Zhao, E. Wijmans, B. Jain, J. Straub, J. Liu, V. Koltun, J. Malik, *et al.*, "Habitat: A platform for embodied ai research," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 9339–9347.

[6] C. D. Singh, R. Kumari, C. Fermüller, N. J. Sanket, and Y. Aloimonos, "Worldgen: A large scale generative simulator," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 9147–9154.

[7] B. Kerbl, G. Kopanas, T. Leimkühler, and G. Drettakis, "3d gaussian splatting for real-time radiance field rendering," *ACM Transactions on Graphics*, vol. 42, no. 4, pp. 1–14, 2023.

[8] Randi Cabezas, Oren Freifeld, Guy Rosman, John W. Fisher II, "Cabezas_Aerial_Reconstructions_via_2014_CVPR_paper.pdf," 2014.

[9] K. Zhang, J. Sun, and N. Snavely, "Leveraging vision reconstruction pipelines for satellite imagery," in *IEEE International Conference on Computer Vision Workshops*, 2019.

[10] J. L. Schönberger and J.-M. Frahm, "Structure-from-motion revisited," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[11] R. A. Beyer, O. Alexandrov, and S. McMichael, "The ames stereo pipeline: Nasa's open source software for deriving and processing terrain data," *Earth and Space Science*, vol. 5, no. 9, pp. 537–548, 2018.

[12] M. Bosch, Z. Kurtz, S. Hagstrom, and M. Brown, "A multiple view stereo benchmark for satellite imagery," in *2016 IEEE Applied Imagery Pattern Recognition Workshop (AIPR)*. IEEE, 2016, pp. 1–9.

[13] P. d'Angelo and G. Kuschk, "Dense multi-view stereo from satellite imagery," in *2012 IEEE International Geoscience and Remote Sensing Symposium*. IEEE, 2012, pp. 6944–6947.

[14] K. Gong and D. Fritsch, "Dsm generation from high resolution multi-view stereo satellite imagery," *Photogrammetric Engineering & Remote Sensing*, vol. 85, no. 5, pp. 379–387, 2019.

[15] S. Chen, X. Wu, M. W. Mueller, and K. Sreenath, "Real-time geo-localization using satellite imagery and topography for unmanned aerial vehicles," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sept. 2021, pp. 2275–2281.

[16] L. M. Downes, D.-K. Kim, T. J. Steiner, and J. P. How, "City-wide Street-to-Satellite image geolocalization of a mobile ground agent," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2022, pp. 11 102–11 108.

[17] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "NeRF: representing scenes as neural radiance fields for view synthesis," *Commun. ACM*, vol. 65, no. 1, pp. 99–106, Dec. 2021.

[18] M. Toschi, R. De Matteo, R. Spezialetti, D. De Gregorio, L. Di Stefano, and S. Salti, "ReLight my NeRF: A dataset for novel view synthesis and relighting of real world objects," Apr. 2023.

[19] C.-H. Lin, W.-C. Ma, A. Torralba, and S. Lucey, "BARF: Bundle-adjusting neural radiance fields," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2021.

[20] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *arXiv preprint arXiv:2201.05989*, 2022.

[21] R. Marí, G. Facciolo, and T. Ehret, "Sat-nerf: Learning multi-view satellite photogrammetry with transient objects and shadow modeling using rpc cameras," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 1311–1321.

[22] H. Wu, A. Graikos, and D. Samaras, "S-VolSDF: Sparse Multi-View stereo regularization of neural implicit surfaces," Mar. 2023.

[23] A. Chen, Z. Xu, F. Zhao, X. Zhang, F. Xiang, J. Yu, and H. Su, "MVSNeRF: Fast generalizable radiance field reconstruction from multi-view stereo," in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2021.

[24] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. Mildenhall, P. P. Srinivasan, J. T. Barron, and H. Kretzschmar, "Block-nerf: Scalable large scene neural view synthesis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8248–8258.

[25] I. Gilitschenski, G. Rosman, A. Gupta, S. Karaman, and D. Rus, "Deep context maps: Agent trajectory prediction using location-specific latent maps," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5097–5104, 2020.

[26] M. Adamkiewicz, T. Chen, A. Caccavale, R. Gardner, P. Culbertson, J. Bohg, and M. Schwager, "Vision-Only robot navigation in a neural radiance world," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 4606–4613, Apr. 2022.

[27] A. Moreau, N. Piasco, D. Tsishkou, B. Stanciulescu, and A. de La Fortelle, "Lens: Localization enhanced by nerf synthesis," in *Conference on Robot Learning*. PMLR, 2022, pp. 1347–1356.

[28] H.-P. Chiu, V. Murali, R. Villamil, G. D. Kessler, S. Samarasekera, and R. Kumar, "Augmented reality driving using semantic geo-registration," in *2018 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*, 2018, pp. 423–430.

[29] I. D. Miller, F. Cladera, T. Smith, C. J. Taylor, and V. Kumar, "Stronger together: Air-ground robotic collaboration using semantics," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 9643–9650, 2022.

[30] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," Apr. 2015.

[31] A. Amini, G. Rosman, S. Karaman, and D. Rus, "Variational end-to-end navigation and localization," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 8958–8964.

[32] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang, *et al.*, "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.

[33] H. Xu, Y. Gao, F. Yu, and T. Darrell, "End-to-end learning of driving models from large-scale video datasets," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2174–2182.

[34] J. Hawke, R. Shen, C. Gurau, S. Sharma, D. Reda, N. Nikolov, P. Mazur, S. Micklethwaite, N. Griffiths, A. Shah, *et al.*, "Urban driving with conditional imitation learning," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 251–257.

[35] A. Byravan, J. Humplik, L. Hasenclever, A. Brussee, F. Nori, T. Haarnoja, B. Moran, S. Bohez, F. Sadeghi, B. Vujatovic, *et al.*, "Nerf2real: Sim2real transfer of vision-guided bipedal motion skills using neural radiance fields," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 9362–9369.

[36] S. Lowry, N. Sünderhauf, P. Newman, J. J. Leonard, D. Cox, P. Corke, and M. J. Milford, "Visual place recognition: A survey," *IEEE Trans. Rob.*, vol. 32, no. 1, pp. 1–19, Feb. 2016.

[37] A. Kendall, M. Grimes, and R. Cipolla, "Posenet: A convolutional network for real-time 6-dof camera relocalization," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2938–2946.

[38] X. Zhang, L. Wang, and Y. Su, "Visual place recognition: A survey from deep learning perspective," *Pattern Recognit.*, vol. 113, p. 107760, May 2021.

[39] R. C. Coulter, "Implementation of the pure pursuit path tracking algorithm," 1992.

[40] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 2019, pp. 8024–8035. [Online]. Available: http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf

[41] 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

[42] N. Buckman, A. Hansen, S. Karaman, and D. Rus, "Evaluating autonomous urban perception and planning in a 1/10th scale minicity," *Sensors*, vol. 22, no. 18, p. 6793, 2022.

[43] "MIT Racecar Platform," @mischttps://mit-racecar.github.io, howpublished = https://mit-racecar.github.io, note = Accessed: 2023-09-15 , accessed: 2023-09-15.